

# Estimating discriminatory power and PD curves when the number of defaults is small

Dirk Tasche, Lloyds Banking Group\*

## Abstract

The intention with this paper is to provide all the estimation concepts and techniques that are needed to implement a two-phases approach to the parametric estimation of probability of default (PD) curves. In the first phase of this approach, a raw PD curve is estimated based on parameters that reflect discriminatory power. In the second phase of the approach, the raw PD curve is calibrated to fit a target unconditional PD. The concepts and techniques presented include a discussion of different definitions of area under the curve (AUC) and accuracy ratio (AR), a simulation study on the performance of confidence interval estimators for AUC, a discussion of the one-parametric approach to the estimation of PD curves by van der Burgt (2008) and alternative approaches, as well as a simulation study on the performance of the presented PD curve estimators. The topics are treated in depth in order to provide the full rationale behind them and to produce results that can be implemented immediately.

## 1 Introduction

In the current economic environment with its particular consequence of rising credit default rates all over the world, at first glance it might not seem very appropriate to look after estimation issues experienced in portfolios with a small number of defaults. However, low default estimation issues can occur quite naturally even in such a situation:

- It is of interest to estimate *instantaneous* discriminatory power of a score function or rating system. “Instantaneous” means that one looks only at the defaults and survivals that occurred in a relatively short time period as one year or less. In typical wholesale portfolios that represent the scope of a rating system the number of borrowers does not exceed 1000. As a consequence, the number of defaults observed within a one year period might well be less than 20.
- Similarly, when estimating forward-looking point-in-time (PIT) conditional probabilities of default per score or rating grade (*PD curve*), it makes sense to construct the estimation sample from observations in a relatively short time period like one or two years in order to capture the instantaneous properties of a potentially rather volatile object. The previous observation on the potentially low number of defaults then applies again.

---

\*The opinions expressed in this paper are those of the author and do not necessarily reflect views of Lloyds Banking Group.

The topics dealt with in this paper are closely related to these two issues. The intention with the paper is to clarify conceptual issues of how to estimate discriminatory power and PD curves, to provide ready-to-use formulas for the related concepts, and to look at low-default related performance issues by means of simulation studies.

The paper is organised as follows:

- **Section 2:** We introduce the concept of a two-phases approach to the calibration of score functions and rating systems and present a simple probabilistic model that is appropriate as a framework to discuss the two phases in a consistent manner. Additionally, we introduce some technical notation for further use within the paper. The focus in this paper will be on the *estimation phase* whose technical details are studied in sections 3, 4, and 5. For the purpose of reference, technical details on the *calibration phase* are provided in appendix A.
- **Section 3:** The estimation of discriminatory power of a score function or a rating system represents an important part of the estimation phase. In particular, when sample sizes are small, it is therefore crucial to have a clear and consistent view on which definition of discriminatory power should be used. This question is discussed in depth in section 3.
- **Section 4:** We replicate simulation studies by Engelmann et al. (2003a,b) with some refinements in order to investigate how accurate confidence interval calculations with different methods for discriminatory power are. In contrast to the studies by Engelmann et al. (2003a,b), the current study is arranged in such a way that the true value of discriminatory power is known. Additionally, even smaller sample sizes are considered and the results are double-checked against exact results from application of the Mann-Whitney test.
- **Section 5:** Van der Burgt (2008) suggested fitting cumulative accuracy profile (CAP) curves with a one-parameter family of exponential curves in order to derive conditional probabilities of default by taking the derivatives of the estimated curves. Van der Burgt believes that this approach is appealing in particular for low default portfolios where the defaulter sample size is small. In section 5, we investigate the suitability of van der Burgt's approach by applying it to a situation where the conditional score distributions are known and, hence, the conditional probabilities of default can be calculated exactly. Performance of van der Burgt's estimator is compared to the performance of the logit estimator and two modifications of it. While it turns out that there is no uniformly best estimator, the results also demonstrate the close relationship between discriminatory power as measured by area under the curve or accuracy ratio and parametric PD curve estimators. A consequence of this, however, is high sensitivity of parametric estimates of PD curves to poor specifications of discriminatory power.
- **Section 6:** Conclusions.

## 2 Basic concepts and notation

Within this paper, we assume that there is a score function or a rating system that can be deployed to inform credit-related decisions. We do not discuss the question of how such score functions or rating systems can be developed. See, e.g., Engelmann and Rauhmeier (2006) for some approaches to the development of rating systems. Instead, in this paper, we look at the

questions of how the power of such score functions or rating systems can be assessed and how probability of default (PD) estimates (*PD curves*) associated with score values or rating grades can be derived.

In sub-section 2.1 we present a general concept for the calibration of score functions and rating systems which is based on separate estimation of discriminatory power and an unconditional probability of default. In sub-section 2.2 a simple probabilistic model is introduced that will help in the derivation of some ideas and formulas needed for implementation of the concept. Moreover, in sub-section 2.3 we recall for further reference some properties of distribution functions and some notation related to such functions.

## 2.1 Estimation phase and calibration phase

In this sub-section, we introduce the concept of a two-phases approach to the calibration of a score function or a rating system: The first phase is the *estimation* phase, the second phase is the *calibration and forecast* phase.

### 2.1.1 Estimation

The aim here is to estimate conditional PDs per score (or grade) and the discriminatory power of the rating system (to be formally defined in sections 3 and 3.2) from a historical sample of scores or rating grades associated with borrowers whose solvency states one period after the scores were observed are known. The composition of the sample is not assumed to be representative of current or future portfolio composition. In particular, the proportion of defaulters and survivors in the sample may differ from proportions of defaulters and survivors expected for the future. The estimated conditional PDs therefore are considered *raw PDs* and have to be *calibrated* before being further used. The estimation sample could be the development sample of the rating system or a validation sample.

In the following, we will write  $x_1, \dots, x_{n_D}$  when talking about a sample of scores or rating grades of defaulted borrowers and  $y_1, \dots, y_{n_N}$  when talking about a sample of surviving borrowers. In both these cases, the solvency state of the borrowers one period after the observation of the scores is known. In contrast, we will write  $s_1, \dots, s_n$  when talking about a sample of scores of borrowers with unknown future solvency state.

### 2.1.2 Calibration and forecast

The aim here is to calibrate the raw PDs from the estimation step in such a way that, on the current portfolio, they are consistent with an unconditional PD that may be different to the unconditional PD of the estimation sample. This calibration exercise is needed because for the borrowers in the current portfolio scores (or rating grades) can be determined but not their future solvency states. Hence direct estimation of conditional PDs with the current portfolio as sample is not possible. We will provide the details of the calibration under the assumption that the *conditional* score distributions (formally defined in (2.3) below) that underlie the estimation sample and the conditional score distributions of the current portfolio are the same. This assumption is reasonable if the estimation sample was constructed not too far back in time or if

the rating system was designed with an intention of creating a through-the-cycle (TTC) rating system.

As mentioned before, the unconditional PDs of estimation sample and current portfolio may be different. This will be the case in particular if a point-in-time (PIT) calibration of the conditional PDs is intended, such that the PDs can be used for forecasting future default rates. But also if a TTC calibration of the PDs is intended (such that no direct forecast of default rates is possible), most of the time the TTC unconditional PD will be different to the realised unconditional PD of the estimation sample. Note that the *unconditional* score distributions of estimation sample and current portfolio can be different but, on principle, are linked together by equation (2.4) from sub-section 2.2.

The question of how to forecast the unconditional PD is not treated in this paper. An example of how PIT estimation of the unconditional PD could be conducted is presented by Engelmann and Porath (2003, section III). Technical details of how the calibration of the conditional PDs can be done are provided in appendix A.

## 2.2 Model and basic properties

Speaking in technical terms, in this paper we study the joint distribution and some estimation aspects of a pair  $(S, Z)$  of real random variables. The variable  $S$  is interpreted as the *credit score* (continuous case) or *rating grade*<sup>1</sup> (discrete case) observed for a solvent borrower at a certain point in time. Hence  $S$  typically takes on values on a continuous scale in some open interval  $I \subset \mathbb{R}$  or on a discrete scale in a finite set  $I = \{1, 2, \dots, k\}$ .

**Convention:** Low values of  $S$  indicate low creditworthiness (“bad”), high values of  $S$  indicate high creditworthiness (“good”).

The variable  $Z$  is the *borrower’s state of solvency* one observation period (usually one year) after the score was observed.  $Z$  takes on values in  $\{0, 1\}$ . The meaning of  $Z = 0$  is “borrower has remained solvent” (solvency or survival),  $Z = 1$  means “borrower has become insolvent” (default). We write  $D$  for the event  $\{Z = 1\}$  and  $N$  for the event  $\{Z = 0\}$ . Hence

$$D \cap N = \{Z = 1\} \cap \{Z = 0\} = \emptyset, \quad D \cup N = \text{whole space.} \quad (2.1)$$

The marginal distribution of the state variable  $Z$  is characterised by the *unconditional probability of default*  $p$  which is defined as

$$p = \mathbb{P}[D] = \mathbb{P}[Z = 1] \in [0, 1]. \quad (2.2)$$

The joint distribution of  $(S, Z)$  then can be specified by the two conditional distributions of  $S$  given the states of  $Z$  or the events  $D$  and  $N$  respectively. In particular, we define the conditional distribution functions

$$\begin{aligned} F_N(s) &= \mathbb{P}[S \leq s | N] = \frac{\mathbb{P}[\{S \leq s\} \cap N]}{1 - p}, \quad s \in I, \\ F_D(s) &= \mathbb{P}[S \leq s | D] = \frac{\mathbb{P}[\{S \leq s\} \cap D]}{p}, \quad s \in I. \end{aligned} \quad (2.3)$$

---

<sup>1</sup>In practice, often a rating system with a small finite number of grades is derived from a score function with values on a continuous scale. This is usually done by mapping score intervals on rating grades. See Tasche (2008, section 3) for a discussion of how such mappings can be defined. Discrete rating systems are preferred by practitioners because manual adjustment of results (overrides) is feasible. Moreover, results by discrete rating systems tend to be more stable over time.

For the sake of an easier notation we denote by  $S_N$  and  $S_D$  random variables with distributions  $\mathbb{P}[S \in \cdot | N]$  and  $\mathbb{P}[S \in \cdot | D]$  respectively. In the literature,  $F_N(s)$  sometimes is called *false alarm rate* while  $F_D(s)$  is called *hit rate*.

By the law of total probability, the distribution function  $F(s) = \mathbb{P}[S \leq s]$  of the marginal (or unconditional) distribution of the score  $S$  can be represented as

$$F(s) = p F_D(s) + (1 - p) F_N(s), \quad \text{all } s. \quad (2.4)$$

$F(s)$  is often called *alarm rate*.

The joint distribution of the pair  $(S, Z)$  of score and borrower's state one period later can also be specified by starting with the unconditional distribution  $\mathbb{P}[S \in \cdot]$  of  $S$  and combining it with the *conditional probability of default*  $\mathbb{P}[D | S] = 1 - \mathbb{P}[N | S]$ . Recall that in general the conditional probability  $\mathbb{P}[D | S] = p_D(S)$  can be characterised<sup>2</sup> by the property (see, e.g. Durrett, 1995, section 4.1)

$$\mathbb{E}[p_D(S) \mathbf{1}_{\{S \in A\}}] = \mathbb{P}[D \cap \{S \in A\}], \quad (2.5)$$

for all Borel sets  $A \subset \mathbb{R}$ . It is well-known (Bayes' formula) that equation (2.5) implies closed-form representations of  $\mathbb{P}[D | S = s] = p_D(s)$  in two important special cases:

- $S$  is a discrete variable, i.e.  $S \in I = \{1, 2, \dots, k\}$ . Then

$$\mathbb{P}[D | S = j] = \frac{p \mathbb{P}[S = j | D]}{p \mathbb{P}[S = j | D] + (1 - p) \mathbb{P}[S = j | N]}, \quad j \in I. \quad (2.6a)$$

- $S$  is a continuous variable with values in an open interval  $I$  such that there are Lebesgue densities  $f_N$  and  $f_D$  of the conditional distribution functions  $F_N$  and  $F_D$  from (2.3). Then

$$\mathbb{P}[D | S = s] = \frac{p f_D(s)}{p f_D(s) + (1 - p) f_N(s)}, \quad s \in I. \quad (2.6b)$$

A closely related consequence of equation (2.5) is the fact that  $p$ ,  $F_N$ , and  $F_D$  can be determined whenever the unconditional score distribution  $F$  and the conditional probabilities of default  $\mathbb{P}[D | S]$  are known. We then obtain

$$p = \mathbb{E}[\mathbb{P}[D | S]] = \begin{cases} \sum_{j=1}^k \mathbb{P}[D | S = j] \mathbb{P}[S = j], & S \text{ discrete} \\ \int_I \mathbb{P}[D | S = s] f(s) ds, & S \text{ continuous with density } f. \end{cases} \quad (2.7a)$$

If  $S$  is a discrete rating variable, we have for  $j \in I$

$$\begin{aligned} \mathbb{P}[S = j | D] &= \mathbb{P}[D | S = j] \mathbb{P}[S = j] / p, \\ \mathbb{P}[S = j | N] &= (1 - \mathbb{P}[D | S = j]) \mathbb{P}[S = j] / (1 - p). \end{aligned} \quad (2.7b)$$

If  $S$  is continuous score variable with density  $f$ , we have for  $s \in I$

$$\begin{aligned} f_D(s) &= \mathbb{P}[D | S = s] f(s) / p, \\ f_N(s) &= (1 - \mathbb{P}[D | S = s]) f(s) / (1 - p). \end{aligned} \quad (2.7c)$$

---

<sup>2</sup>We define the indicator function  $\mathbf{1}_M$  of a set  $M$  by  $\mathbf{1}_M(m) = \begin{cases} 1, & m \in M, \\ 0, & m \notin M. \end{cases}$

### 2.3 Notation for distribution functions

At some points in this paper we will need to handle distribution functions and their inverse functions. For further reference we list in this subsection the necessary notation and some properties of such functions:

- A (real) distribution function  $G$  is an increasing and right-continuous function  $\mathbb{R} \rightarrow [0, 1]$  with  $\lim_{x \rightarrow -\infty} G(x) = 0$  and  $\lim_{x \rightarrow \infty} G(x) = 1$ .
- Any real random variable  $X$  defines a distribution function  $G = G_X$  by  $G(x) = \mathbb{P}[X \leq x]$ .
- **Convention:**  $G(-\infty) = 0$  and  $G(\infty) = 1$ .
- Denote by  $G(\cdot - 0)$  the left-continuous version of the distribution function  $G$ . Then  $G(\cdot - 0) \leq G$  and  $G(x - 0) = G(x)$  for all  $x$  but countably many  $x \in \mathbb{R}$  because  $G$  is non-decreasing.
- For any distribution function  $G$ , the function  $G^{-1}$  is its *generalised inverse* or *quantile function*, i.e.

$$G^{-1}(u) = \inf\{x \in \mathbb{R} : G(x) \geq u\}, \quad u \in [0, 1]. \quad (2.8a)$$

In particular, we obtain

$$-\infty = G^{-1}(0) < G^{-1}(1) \leq \infty. \quad (2.8b)$$

- Denote by  $\varphi(s)$  the standard normal density and by  $\Phi(s)$  the standard normal distribution function.

## 3 Discriminatory power: Theory

Hand (1997, section 8.1) described ROC curves as follows: “Often the two degrees of freedom [i.e. the two error types associated with binary classification] are presented simultaneously for a range of possible classification thresholds for the classifier in a *receiver operating characteristic (ROC) curve*. This is done by plotting true positive rate (sensitivity) on the vertical axis against false positive rate (1 - specificity) on the horizontal axis.”

Translated into the notation introduced in section 2, for a fixed score value  $s$  seen as threshold the true positive rate is the hit rate  $F_D(s)$  while the false positive rate is the false alarm rate  $F_N(s)$ . In these terms, CAP (Cumulative Accuracy Profile) curves (not mentioned by Hand, 1997) can be described as a plot of the hit rates against the alarm rates across a range of classification thresholds. If all possible thresholds are to be considered, these descriptions formally can be expressed in the following terms.

**Definition 3.1 (ROC and CAP)** *Denote by  $F_N$  the distribution function  $F_N(s) = \mathbb{P}[S_N \leq s]$  of the scores conditional on the event “borrower survives”, by  $F_D$  the distribution function  $F_D(s) = \mathbb{P}[S_D \leq s]$  of the scores conditional on the event “borrower defaults”, and by  $F$  the unconditional distribution function  $F(s) = \mathbb{P}[S \leq s]$  of the scores.*

*The Receiver Operating Characteristic (ROC) of the score function then is defined as the graph of the following set gROC (“g” for graph) of points in the unit square:*

$$\text{gROC} = \{(F_N(s), F_D(s)) : s \in \mathbb{R} \cup \{\pm\infty\}\}. \quad (3.1a)$$

The Cumulative Accuracy Profile (AUC) of the score function is defined as the graph of the following set gCAP of points in the unit square:

$$\text{gCAP} = \{(F(s), F_D(s)) : s \in \mathbb{R} \cup \{\pm\infty\}\}. \quad (3.1b)$$

Actually the point sets gROC and gCAP can be quite irregular (e.g. if one of the involved distribution functions has an infinite number of discontinuities and the set of discontinuities is dense in  $\mathbb{R}$ ). In such a case it would be physically impossible to plot on paper a precise graph of the point set. In most parts of the following, therefore, we will focus on three more regular special cases which are of relevance for theory and practice:

- 1)  $F$ ,  $F_N$ , and  $F_D$  are smooth, i.e. at least continuous. This is usually a reasonable assumption when the score function takes on values on a continuous scale.
- 2) The distributions of  $S$ ,  $S_N$  and  $S_D$  are concentrated on a finite number of points. This is the case when the score function is a rating system with a finite number (e.g. seven or seventeen as in case of S & P, Moody's, or Fitch ratings) of grades.
- 3)  $F$ ,  $F_N$ , and  $F_D$  are empirical distribution functions associated to finite samples of scores on a continuous scale. This is naturally the case when the performance of a score function is analysed on the basis of non-parametric estimates.

In the smooth situation of 1) the sets gROC and gCAP are compact and connected such that there is no ambiguity left of how to draw a graph that – together with the x-axis and the vertical line through  $x = 1$  – encloses a region of finite area. In situations 2) and 3), however, the sets gROC and gCAP consist of a finite number of isolated points and hence are unconnected. While this, in a certain sense, even facilitates the drawing of the graphs, the results nonetheless will be unsatisfactory when it comes to a comparison of the discriminatory power of score functions or rating systems. Usually, therefore, in such cases a certain degree of interpolation will be applied to the points of the sets gROC and gCAP in order to facilitate their visual comparison. We will discuss in section 3.2 the question of how to do best the interpolation to satisfy some properties that are desirable from a statistical point of view.

Before, however, in section 3.1 we have a closer look on the properties of ROC graphs in smooth contexts. These properties then will be used as a kind of yardstick to assess the appropriateness of interpolation approaches to the discontinuous case in section 3.2.

### 3.1 Continuous score distributions

In this subsection, we will work most of the time on the basis of one of the following two assumptions.

**Assumption N:** The distribution of the score  $S_N$  conditional on the borrower's survival is continuous, i.e.

$$P[S_N = s] = 0 \quad \text{for all } s. \quad (3.2)$$

**Assumption S:** The unconditional distribution of the score  $S$  is continuous (and hence by (2.4) so are the distributions of  $S_N$  and  $S_D$ ), i.e.

$$P[S = s] = 0 \quad \text{for all } s. \quad (3.3)$$

Additionally, the following technical assumption is sometimes useful.

**Assumption:**

$$F_D^{-1}(1) \leq F_N^{-1}(1). \quad (3.4)$$

This is equivalent to requiring that the essential supremum of  $S_D$  is not greater than the essential supremum of  $S_N$ . Such a requirement seems natural under the assumption that low score values indicate low creditworthiness (“bad”) and high score values indicate high creditworthiness (“good”).

As an immediate consequence of these assumptions we obtain representations of the ROC and CAP sets (3.1a) and (3.1b) that are more convenient for calculations.

**Theorem 3.2 (Standard parametrisations of ROC and CAP)**

With the notation of definition 3.1 define the functions ROC and CAP by

$$\text{ROC}(u) = F_D(F_N^{-1}(u)), \quad u \in [0, 1], \quad (3.5a)$$

$$\text{CAP}(u) = F_D(F^{-1}(u)) = F_D((p F_D(\cdot) + (1 - p) F_N(\cdot))^{-1}(u)), \quad u \in [0, 1]. \quad (3.5b)$$

For (3.5b), assume  $p > 0$  (otherwise ROC and CAP coincide). Under (3.2) (assumption N) then we have

$$\{(u, \text{ROC}(u)) : u \in [0, 1]\} \subset \text{gROC}. \quad (3.5c)$$

If under (3.2) (assumption N), moreover, the distribution of  $S_D$  is absolutely continuous with respect to the distribution of  $S_N$  (i.e.  $\mathbb{P}[S_N \in A] = 0 \Rightarrow \mathbb{P}[S_D \in A] = 0$ ), then<sup>3</sup> “ $\subset$ ” applies also to (3.5c):

$$\{(u, \text{ROC}(u)) : u \in [0, 1]\} = \text{gROC}. \quad (3.5d)$$

Equation (3.3) (assumption S) implies

$$\{(u, \text{CAP}(u)) : u \in [0, 1]\} = \text{gCAP}. \quad (3.5e)$$

**Proof.** Note that (2.8b) implies in general

$$0 = \text{CAP}(0) = \text{ROC}(0). \quad (3.6a)$$

For  $p > 0$ , we have  $\{s : p F_D(s) + (1 - p) F_N(s) \geq 1\} \subset \{s : F_D(s) \geq 1\}$  and hence

$$\begin{aligned} \text{CAP}(1) &= F_D((p F_D(\cdot) + (1 - p) F_N(\cdot))^{-1}(1)) \\ &\geq F_D(F_D^{-1}(1)) \\ &\geq 1 \\ \Rightarrow \text{CAP}(1) &= 1. \end{aligned} \quad (3.6b)$$

Additionally, if (3.4) holds – which is implied by the absolute continuity assumption – we obtain

$$\begin{aligned} \text{ROC}(1) &= F_D(F_N^{-1}(1)) \\ &\geq F_D(F_D^{-1}(1)) \\ &\geq 1 \\ \Rightarrow \text{ROC}(1) &= 1. \end{aligned} \quad (3.6c)$$

---

<sup>3</sup>The absolute continuity requirement implies that  $F_D$  is constant on the intervals on which  $F_N$  is constant.



Now, by (3.2) (assumption N) we have  $F_N(F_N^{-1}(u)) = u$  and by (3.3) (assumption S) we have  $F(F^{-1}(u)) = u$  (see van der Vaart, 1998, section 21.1). This implies (3.5c) and “ $\subset$ ” in (3.5e). Assume that distribution of  $S_D$  is absolutely continuous with respect to the distribution of  $S_N$ . For  $s \in \mathbb{R}$  let  $s_0 = F_N^{-1}(F_N(s))$ . By continuity of  $F_N$  then we have  $F_N(s_0) = F_N(s)$ , and by absolute continuity of  $S_D$  with respect to  $S_N$  we also have  $F_D(s_0) = F_D(s)$ . This implies “ $=$ ” in (3.5c) because on the one hand

$$(F_N(s), F_D(s)) = (F_N(s), F_D(s_0)) = (F_N(s), \text{ROC}(F_N(s))) \in \{(u, \text{ROC}(u)) : u \in [0, 1]\},$$

and on the other hand for  $s = \pm\infty$  we can apply (3.6a), (3.6b), and (3.6c).

The “ $=$ ” in (3.5e) follows from the fact that  $S_D$  by (2.4) is always absolutely continuous with respect to  $S$ .  $\square$

### Remark 3.3

A closer analysis of the proof of theorem 3.2 shows that a non-empty difference between the left-hand and the right-hand sides of (3.5c) can occur only if there are non-empty intervals on which the value of  $F_N$  is constant. To each such interval on which  $F_D$  is not constant there is corresponding piece of a vertical line in the set  $\text{gROC}$  that has no counterpart in the graph of the function  $\text{ROC}(u)$ . Note, however, that these missing pieces are not relevant with respect to the area below the ROC curve because this area is still well-defined when all vertical pieces are removed from  $\text{gROC}$ . In this sense, in theorem 3.2 the absolute continuity requirement and equation (3.5d) are only of secondary importance.

In view of theorem 3.2 and remark 3.3, we can regard ROC and CAP curves as graphs for functions (3.5a) and (3.5b) respectively, as long as (3.2) and (3.3) apply. This provides a convenient way to dealing analytically with ROC and CAP curves. In section 3.2 we will revisit the question of how to conveniently parametrize the point sets (3.1a) and (3.1b) in the case of score distributions with discontinuities.

In this section, we continue by looking closer at some well-known properties of ROC and CAP curves. In non-technical terms the following proposition 3.4 states: The diagonal line is the ROC and CAP curve of powerless rating systems (or score functions). For a perfect score function, the ROC curve is essentially the horizontal line at level 1 while the CAP curve is made up by the straight line  $u \mapsto u/p, u < p$  and the horizontal line at level 1.

**Proposition 3.4** *Under (3.2) (assumption N), in case of a powerless classification system (i.e.  $F_D = F_N$ ) we have*

$$\text{ROC}(u) = u = \text{CAP}(u), \quad u \in [0, 1]. \quad (3.7a)$$

*In case of a perfect classification system<sup>4</sup> (i.e. there is a score value  $s_0$  such that  $F_D(s_0) = 1, F_N(s_0) = 0$ ) we obtain without continuity assumption that*

$$\text{ROC}(u) = \begin{cases} 0, & u = 0, \\ 1, & 0 < u \leq 1, \end{cases} \quad (3.7b)$$

*and, if  $p > 0$  and  $F_D$  is continuous,*

$$\text{CAP}(u) = \begin{cases} u/p, & 0 \leq u < p, \\ 1, & p \leq u \leq 1. \end{cases} \quad (3.7c)$$

---

<sup>4</sup>Note that in case of a perfect classification system the distribution of  $S_D$  is not absolutely continuous with respect to the distribution of  $S_N$  as it would be required for (3.5d) to obtain.

**Proof.** For (3.7a), we have to show that

$$F_D(F_D^{-1}(u)) = u, \quad u \in [0, 1]. \quad (3.8)$$

This follows from the continuity assumption (3.2) (see van der Vaart, 1998, section 21.1).

On (3.7b) and (3.7c): Observe that  $F_D(s_0) = 1, F_N(s_0) = 0$  for some  $s_0$  implies (3.4). By (3.6a), (3.6b) and (3.6c), therefore, we only need to consider the case  $0 < u < 1$ . For  $u > p$  we obtain

$$\begin{aligned} F(s_0) &= p F_D(s_0) + (1 - p) F_N(s_0) = p \\ \Rightarrow F^{-1}(u) &\geq s_0 \\ \Rightarrow F_D(F^{-1}(u)) &= 1. \end{aligned} \quad (3.9)$$

This implies (3.7b) (with  $p = 0$ ), in particular, and (3.7c) for  $u > p$ . For  $u < p$ , equation (3.9) implies  $F^{-1}(u) < s_0$ . By left continuity of  $F^{-1}$ , we additionally obtain  $F^{-1}(u) \leq s_0$  for  $u \leq p$ . But

$$F(s) = p F_D(s) + (1 - p) F_N(s) = p F_D(s), \quad s \leq s_0.$$

Hence for  $u \leq p$

$$\begin{aligned} F^{-1}(u) &= \inf\{s : p F_D(s) \geq u\} = F_D^{-1}(u/p) \\ \Rightarrow F_D(F^{-1}(u)) &= F_D(F_D^{-1}(u/p)) = u/p. \end{aligned}$$

The last equality follows from the assumed continuity of  $F_D$ . □

By theorem 3.2, in the continuous case (3.2) and (3.3), the common notions of AUC (area under the curve) and AR (accuracy ratio) can be defined in terms of integrals of the ROC and CAP functions (3.5a) and (3.5b). Recall that the accuracy ratio commonly is described in terms like these: “The quality of a rating system is measured by the accuracy ratio AR. It is defined as the ratio of the area between the CAP of the rating model being validated and the CAP of the random model [= powerless model], and the area between the CAP of the perfect rating model and the CAP of the random model” (Engelmann et al., 2003a, page 82).

**Definition 3.5 (Area under the curve and accuracy ratio)**

For the function ROC given by (3.5a) we define the area under the curve AUC by

$$\text{AUC} = \int_0^1 \text{ROC}(u) du. \quad (3.10a)$$

For the function CAP given by (3.5b) we define the accuracy ratio AR by

$$\text{AR} = \frac{\int_0^1 \text{CAP}(u) - u du}{1 - p/2 - 1/2} = \frac{2 \int_0^1 \text{CAP}(u) du - 1}{1 - p}. \quad (3.10b)$$

In the continuous case (3.3) (assumption S) AUC and AR are identical up to a constant linear transformation, as shown by the following proposition.

**Proposition 3.6 (AUC and AR in the continuous case)**

If the distribution of the score function conditional on default is continuous then

$$\text{AR} = 2 \text{AUC} - 1.$$

**Proof.** Denote by  $S'_D$  a random variable with the same distribution  $F_D$  as  $S_D$  but independent of  $S_D$ . Let  $S_N$  be independent of  $S_D$ . Observe that  $F_N^{-1}(U)$  and  $F_D^{-1}(U)$  have the same distribution as  $S_N$  and  $S_D$  if  $U$  is uniformly distributed on  $(0, 1)$ . By the definition (3.10b) of AR and Fubini's theorem, therefore we obtain

$$\begin{aligned} \text{AR} &= \frac{2}{1-p} (p \text{P}[S_D \leq S'_D] + (1-p) \text{P}[S_D \leq S_N] - 1/2) \\ &= 2 \text{P}[S_D \leq S_N] - 1 \\ &= 2 \text{AUC} - 1. \end{aligned} \tag{3.11}$$

In this calculation, the fact has been used that  $1/2 = \text{P}[S_D \leq S'_D]$  because the distribution of  $S_D$  is assumed to be continuous.  $\square$

As the ROC curve does not depend on the proportion  $p$  of defaulters in the population, proposition 3.6 in particular shows that AR does not depend on  $p$  either. The following corollary is an easy consequence of propositions 3.4 and 3.6. It identifies the extreme cases for classification systems. A classification system is considered poor if its AUC and AR are close to AUC and AR of a powerless system. It is considered powerful if its AUC and AR are close to AUC and AR of a perfect system.

**Corollary 3.7** *Under (3.2) (assumption N), in case of a powerless classification system (i.e.  $F_D = F_N$ ) we have*

$$\begin{aligned} \text{AUC} &= 1/2, \\ \text{AR} &= 0. \end{aligned} \tag{3.12a}$$

*In case of a perfect classification system (i.e. there is a score value  $s_0$  such that  $F_D(s_0) = 1, F_N(s_0) = 0$ ) we obtain if the distribution of the scores conditional on default is continuous*

$$\begin{aligned} \text{AUC} &= 1, \\ \text{AR} &= 1. \end{aligned} \tag{3.12b}$$

Relation (3.12a) can obtain also in situations where  $F_N \neq F_D$ . For instance, Clavero Rasero (2006, proposition 2.6) proved that (3.12a) applies in general when  $F_N$  and  $F_D$  have densities that are both symmetric with respect to the same point.

### 3.1.1 Example: Normally distributed scores

Assume that the score distributions conditional on default and survival, respectively, are normal:

$$S_D \sim \mathcal{N}(\mu_D, \sigma_D^2), \quad S_N \sim \mathcal{N}(\mu_N, \sigma_N^2). \tag{3.13}$$

Formulas for ROC in the sense of (3.5a) and AUC in the sense of (3.10a) then easily are derived:

$$\begin{aligned} \text{ROC}(u) &= \Phi\left(\frac{\sigma_N \Phi^{-1}(u) + \mu_N - \mu_D}{\sigma_D}\right), \quad u \in [0, 1] \\ \text{AUC} &= \Phi\left(\frac{\mu_N - \mu_D}{\sqrt{\sigma_N^2 + \sigma_D^2}}\right). \end{aligned} \tag{3.14}$$

Note that (3.14) gives a closed form of AUC where Satchell and Xia (2008) provided a formula involving integration. See figure 1 for an illustration of (3.13) and (3.14).

The unconditional score distribution  $F$  can be derived from (2.4). Under (3.13), however, for  $p \notin \{0, 1\}$ ,  $F$  is not a normal distribution function. Its inverse function  $F^{-1}$  can be evaluated numerically, but no closed-form representation is known. For plots of the CAP curve, therefore it is more efficient to make use of representation (3.1b). The value of AR can be derived from the value of AUC by proposition 3.6.

### 3.1.2 Example: Density estimation with normal kernel

Assume that there are samples  $x_1, \dots, x_{n_D}$  of scores of defaulted borrowers and  $y_1, \dots, y_{n_N}$  of surviving borrowers. If the scores take on values on a continuous scale, it makes sense to try and estimate densities of the defaulters' scores and survivors' scores, respectively. We consider here kernel estimation with a normal kernel as estimation approach (see, e.g. Pagan and Ullah, 1999, chapter 2). The resulting density estimates then are

$$\begin{aligned}\widehat{f}_D(s) &= (n_D h_D)^{-1} \sum_{i=1}^{n_D} \varphi\left(\frac{s - x_i}{h_D}\right), \\ \widehat{f}_N(s) &= (n_N h_N)^{-1} \sum_{i=1}^{n_N} \varphi\left(\frac{s - y_i}{h_N}\right),\end{aligned}\tag{3.15}$$

where  $h_D, h_N > 0$  denote appropriately selected *bandwidths*. Silverman's rule of thumb (see, e.g. Pagan and Ullah, 1999, equation (2.50)) often yields reasonable results:

$$h = 1.06 \widehat{\sigma} T^{-1/5},\tag{3.16}$$

where  $\widehat{\sigma}$  denotes the standard deviation of the sample  $x_1, \dots, x_{n_D}$  or  $y_1, \dots, y_{n_N}$ , respectively. Equation (3.15) immediately implies the following formulas for the corresponding estimated distribution functions:

$$\begin{aligned}\widehat{F}_D(s) &= (n_D)^{-1} \sum_{i=1}^{n_D} \Phi\left(\frac{s - x_i}{h_D}\right), \\ \widehat{F}_N(s) &= (n_N)^{-1} \sum_{i=1}^{n_N} \Phi\left(\frac{s - y_i}{h_N}\right).\end{aligned}\tag{3.17}$$

ROC and CAP curves then can be drawn efficiently by taking recourse to (3.1a) and (3.1b). An estimate of AUC (and then by proposition 3.6 of AR) is given by a generalisation of (3.14):

$$\widehat{AUC} = (n_D n_N)^{-1} \sum_{i=1}^{n_D} \sum_{j=1}^{n_N} \Phi\left(\frac{y_j - x_i}{\sqrt{h_N^2 + h_D^2}}\right).\tag{3.18}$$

See figure 2 for illustration.

#### Remark 3.8 (Bias of kernel-based AUC-estimator)

Assume that the samples  $x_1, \dots, x_{n_D}$  of scores of defaulted borrowers and  $y_1, \dots, y_{n_N}$  of scores

of surviving borrowers are samples from normally distributed score functions as in (3.13). Then the expected value of the AUC-estimator  $\widehat{AUC}$  from (3.18) can be calculated as follows:

$$\mathbb{E}[\widehat{AUC}] = \Phi\left(\frac{\mu_N - \mu_D}{\sqrt{h_N^2 + h_D^2 + \sigma_N^2 + \sigma_D^2}}\right).$$

Hence by (3.14), the following observations apply:

$$\begin{aligned} |\mathbb{E}[\widehat{AUC}] - 1/2| &\leq |\text{AUC} - 1/2| \\ \text{sign}(\mathbb{E}[\widehat{AUC}] - 1/2) &= \text{sign}(\text{AUC} - 1/2) \\ \mu_D = \mu_N &\Leftrightarrow \mathbb{E}[\widehat{AUC}] = \text{AUC} \\ \mu_D = \mu_N &\Leftrightarrow \text{AUC} = 1/2. \end{aligned}$$

In particular, in case  $\mu_N > \mu_D$  the estimator  $\widehat{AUC}$  on average underestimates the area under the curve while in case  $\mu_N < \mu_D$  the area under the curve is overestimated by  $\widehat{AUC}$ .

To account for the potential bias of the AUC estimates by (3.18) as observed in remark 3.8, in section 4 we will apply linear transformations to the density estimates (3.15). These linear transformations make sure that the means and variances of the estimated densities exactly match the empirical means and variances of the samples  $x_1, \dots, x_{n_D}$  and  $y_1, \dots, y_{n_N}$  respectively (Davison and Hinkley, 1997, section 3.4). Define

$$b_D = \sqrt{\frac{1/n \sum_{i=1}^{n_D} x_i^2 - (1/n \sum_{i=1}^{n_D} x_i)^2}{h_D^2 + 1/n \sum_{i=1}^{n_D} x_i^2 - (1/n \sum_{i=1}^{n_D} x_i)^2}}, \quad a_D = \frac{1 - b_D}{n} \sum_{i=1}^{n_D} x_i, \quad (3.19a)$$

$$b_N = \sqrt{\frac{1/n \sum_{j=1}^{n_N} y_j^2 - (1/n \sum_{j=1}^{n_N} y_j)^2}{h_N^2 + 1/n \sum_{j=1}^{n_N} y_j^2 - (1/n \sum_{j=1}^{n_N} y_j)^2}}, \quad a_N = \frac{1 - b_N}{n} \sum_{j=1}^{n_N} y_j. \quad (3.19b)$$

Replace then in equations (3.15), (3.17), and (3.18)

$$\begin{aligned} x_i &\text{ by } a_D + b_D x_i \quad \text{and} \quad h_D \text{ by } b_D h_D, \\ y_j &\text{ by } a_N + b_N y_j \quad \text{and} \quad h_N \text{ by } b_N h_N, \end{aligned} \quad (3.19c)$$

to reduce the bias from an application of (3.18) for AUC estimation. If, for instance, in the right-hand panel of figure 2 the estimated ROC curve is based on the transformed samples according to (3.19c), the resulting estimate of AUC is 71.2%. Thus, at least in this example, the “transformed” AUC estimate is closer to the true value of 71.6% than the estimate based on estimated densities without adjustments for mean and variance.

## 3.2 Discontinuous score distributions

We have seen that in the case of continuous score distributions as considered in section 3.1 there are standard representations of ROC and CAP curves (theorem 3.2) that can be conveniently deployed to formally define the area under the curve (AUC) and the accuracy ratio (AR) and to investigate some of their properties. In this section, we will see that in a more general setting

the use of the curve representations (3.5a) and (3.5b) can have counter-intuitive implications. We then will look at modifications of (3.5a) and (3.5b) that avoid such implications and show that these modifications are compatible with common interpolation approaches to the ROC and CAP graphs as given by (3.1a) and (3.1b). We will do so primarily with a view on the settings described in items 2) and 3) at the beginning of section 3. For the sake of reference, the following two examples describe these settings in more detail.

**Example 3.9 (Rating distributions)**

Consider a rating system with grades  $1, 2, \dots, n$  where  $n$  stands for highest creditworthiness. The random variable  $R$  which expresses a borrower's rating grade then is purely discontinuous because

$$\begin{aligned} P[R = k] &\geq 0, \quad k \in \{1, 2, \dots, n\} \\ P[R \notin \{1, 2, \dots, n\}] &= 0. \end{aligned}$$

See the upper panel of figure 3 for illustration. As in the case of score functions  $S$ , we write  $R_D$  when considering  $R$  on the sub-population of defaulters and  $R_N$  when considering  $R$  on the sub-population of survivors.

**Example 3.10 (Sample-based empirical distributions)**

Assume – as in section 3.1.2 – that there are samples  $x_1, \dots, x_{n_D}$  of scores of defaulted borrowers and  $y_1, \dots, y_{n_N}$  of surviving borrowers. If there is no reason to believe that the samples were generated from continuous score distributions, or if sample sizes are so large that kernel estimation becomes numerically inefficient, one might prefer to work with the empirical distributions of  $S_D$  and  $S_N$  as inferred from  $x_1, \dots, x_{n_D}$  and  $y_1, \dots, y_{n_N}$ , respectively:

For  $w, z \in \mathbb{R}$  let

$$\delta_w(z) = \begin{cases} 1, & z \leq w \\ 0, & z > w. \end{cases}$$

For  $w \in \mathbb{R}$  define the empirical distribution function for the sample  $z_1, \dots, z_n$  by

$$\delta_w(z_1, \dots, z_n) = 1/n \sum_{i=1}^n \delta_w(z_i). \tag{3.20a}$$

For  $w, z \in \mathbb{R}$  let

$$\delta_w^*(z) = \begin{cases} 1, & z < w \\ 1/2, & z = w \\ 0, & z > w. \end{cases}$$

For  $w \in \mathbb{R}$  define the modified empirical distribution function for the sample  $z_1, \dots, z_n$  by

$$\delta_w^*(z_1, \dots, z_n) = 1/n \sum_{i=1}^n \delta_w^*(z_i). \tag{3.20b}$$

Of course, there is some overlap between examples 3.9 and 3.10. The samples in example 3.10 could have been generated from rating distributions as described in example 3.9 (see lower panel of figure 3 for illustration). Then example 3.10 just would be a special case of example 3.9. The more interesting case in example 3.10 therefore is the case where  $\{x_1, \dots, x_{n_D}\} \cap \{y_1, \dots, y_{n_N}\} =$

$\emptyset$ . This will occur with probability 1 when the two sub-samples are generated from continuous score distributions.

**Some consequences of discontinuity:**

- In the settings of examples 3.9 and 3.10 the CAP and ROC graphs as defined by (3.1b) and (3.1a) consist of finitely many points.
- CAP and ROC functions as defined by (3.5b) and (3.5a) are piecewise constant for rating grade variables  $R$  as in example 3.9 and empirical distribution functions as in example 3.10. See left panel of figure 4 for illustration.
- Proposition 3.4 does not apply. In particular, the graphs of CAP and ROC functions as defined by (3.5b) and (3.5a) for powerless score functions with discontinuities are not identical with the diagonal line. See left panel of figure 4 for illustration.
- Let  $S$  be a random variable with a distribution that is concentrated on finitely many points as in example 3.9 or 3.10. Let  $S'$  be a random variable with the same distribution as  $S$  but independent of  $S$ . Then we have

$$P[S = S'] > 0. \tag{3.21}$$

**3.2.1 Observations on the general case**

In this section, we first look at what happens with corollary 3.7 if no continuity assumption obtains.

**Proposition 3.11 (AUC and AR in the general case)**

Define AUC and AR by (3.10a) and (3.10b), respectively, with ROC and CAP as given in (3.5a) and (3.5b). Let  $S_D$  and  $S_N$  denote independent random variables with distribution functions  $F_D$  (score distribution conditional on default) and  $F_N$  (score distribution conditional on survival). Assume that  $S'_D$  is an independent copy of  $S_D$ . Then

$$\begin{aligned} \text{AUC} &= P[S_D \leq S_N], \\ \text{AR} &= 2P[S_D \leq S_N] - 1 + \frac{p}{1-p} P[S_D = S'_D]. \end{aligned}$$

**Proof.** The equation for AUC follows from application of Fubini’s theorem to the right-hand side of (3.10a). Observe that in general

$$2P[S_D \leq S'_D] = 1 + P[S_D = S'_D] \tag{3.22a}$$

and therefore

$$P[S_D \leq S'_D] - 1/2 = P[S_D = S'_D]/2. \tag{3.22b}$$

Inserting this last identity into (3.11) yields the equation for AR. □

**Corollary 3.12** Define AUC and AR by (3.10a) and (3.10b), respectively, with ROC and CAP as given in (3.5a) and (3.5b). Let  $S_D$  and  $S_{D'}$  denote independent random variables with distribution function  $F_D$ . In case of a powerless classification system (i.e.  $F_D = F_N$ ) we then have

$$\begin{aligned} \text{AUC} &= 1/2 + \text{P}[S_D = S'_D]/2, \\ \text{AR} &= \frac{\text{P}[S_D = S'_D]}{1 - p}. \end{aligned} \quad (3.23)$$

In case of a perfect classification system (i.e. there is a score value  $s_0$  such that  $F_D(s_0) = 1, F_N(s_0) = 0$ ) we have

$$\text{AUC} = 1 \quad (3.24a)$$

and, if  $p > 0$ ,

$$\text{AR} = 1 + \frac{p}{1 - p} \text{P}[S_D = S'_D]. \quad (3.24b)$$

When corollary 3.12 is compared to corollary 3.7, it becomes clear that definitions (3.5a) and (3.5b) are unsatisfactory when it comes to calculate AUC and AR for powerless or perfect score functions with potential discontinuities. In particular, AUC and AR of powerless score functions then will not equal any longer 50% and 0, respectively. AR of a perfect score function can even be greater than 100% when calculated for a score function with discontinuities.

Definitions (3.5a) and (3.5b) of ROC and CAP curves, however, can be modified in a way such that proposition 3.6 and corollary 3.7 obtain without the assumption that the score function is continuous.

**Definition 3.13 (Modified ROC and CAP functions)**

Denote by  $F_N$  and  $F_D$  the distribution functions of the survivor scores and the defaulter scores respectively. Let  $S_D$  be a random variable with distribution function  $F_D$ . The Modified Receiver Operating Characteristic function  $\text{ROC}^*(u)$  then is defined by

$$\text{ROC}^*(u) = \text{P}[S_D < F_N^{-1}(u)] + \text{P}[S_D = F_N^{-1}(u)]/2, \quad u \in [0, 1]. \quad (3.25a)$$

With  $F$  denoting the unconditional distribution function of the scores, the Modified Cumulative Accuracy Profile function  $\text{CAP}^*(u)$  is defined by

$$\text{CAP}^*(u) = \text{P}[S_D < F^{-1}(u)] + \text{P}[S_D = F^{-1}(u)]/2, \quad u \in [0, 1]. \quad (3.25b)$$

In general, we have

$$\text{ROC}^*(u) \leq \text{ROC}(u) \quad \text{and} \quad \text{CAP}^*(u) \leq \text{CAP}(u), \quad u \in [0, 1].$$

Compare the two panels of figure 4 for illustration. If, however, the distribution function  $F_D$  of the defaulter scores is continuous, (3.25a) and (3.5a) are equivalent, and so are (3.25b) and (3.5b) because

$$\begin{aligned} \text{ROC}(u) &= \text{P}[S_D < F_N^{-1}(u)] + \text{P}[S_D = F_N^{-1}(u)], \\ \text{CAP}(u) &= \text{P}[S_D < F^{-1}(u)] + \text{P}[S_D = F^{-1}(u)]. \end{aligned} \quad (3.26)$$



The following modified definitions of AUC and AR obviously coincide with the unmodified concepts of AUC and AR from definition 3.5 when the underlying score distributions are continuous.

**Definition 3.14 (Modified area under the curve and modified accuracy ratio)**

For the function  $\text{ROC}^*$  given by (3.25a) we define the modified area under the curve  $\text{AUC}^*$  by

$$\text{AUC}^* = \int_0^1 \text{ROC}^*(u) du. \quad (3.27a)$$

For the function  $\text{CAP}^*$  given by (3.25b) we define the modified accuracy ratio  $\text{AR}^*$  by

$$\text{AR}^* = \frac{2}{1-p} \left( \int_0^1 \text{CAP}^*(u) du - 1/2 \right). \quad (3.27b)$$

Clearly, we have  $\text{AUC}^* \leq \text{AUC}$  and  $\text{AR}^* \leq \text{AR}$ . The advantage of definition 3.14 compared to definition 3.5 is that it gives us versions of proposition 3.6 and corollary 3.7 that obtain without any continuity requirements on the score distributions.

**Proposition 3.15** Define  $\text{AUC}^*$  and  $\text{AR}^*$  by (3.27a) and (3.27b), respectively, with  $\text{ROC}^*$  and  $\text{CAP}^*$  as given in (3.25a) and (3.25b). Let  $S_D$  and  $S_N$  denote independent random variables that have the distribution of the scores conditional on default and on survival respectively. Then we obtain

$$\text{AUC}^* = \text{P}[S_D < S_N] + \text{P}[S_D = S_N]/2, \quad (3.28a)$$

$$\text{AR}^* = 2 \text{P}[S_D < S_N] + \text{P}[S_D = S_N] - 1 = \text{P}[S_D < S_N] - \text{P}[S_D > S_N]. \quad (3.28b)$$

In particular,  $\text{AR}^* = 2 \text{AUC}^* - 1$  holds.

**Proof.** By application of Fubini's theorem, obvious from the definitions of  $\text{AUC}^*$  and  $\text{AR}^*$ .  $\square$

Note that (3.28a) by some authors (e.g. Newson, 2001, equation (12)) is used as definition of the area under the ROC curve.

**Corollary 3.16** In case of a powerless classification system (i.e.  $F_D = F_N$ ) we have

$$\begin{aligned} \text{AUC}^* &= 1/2, \\ \text{AR}^* &= 0. \end{aligned} \quad (3.29)$$

In case of a perfect classification system (i.e. there is a score value  $s_0$  such that  $F_D(s_0) = 1, F_N(s_0) = 0$ ) we have

$$\text{AUC}^* = 1 \quad (3.30a)$$

and, if  $p > 0$ ,

$$\text{AR}^* = 1. \quad (3.30b)$$

Corollary 3.16 gives a clear indication that for general score distributions definition 3.14 should be preferred to definition 3.5. For the latter definition leads to the results from corollary 3.12 that are counter-intuitive in case of discontinuous score distributions. In section 3.2.2, we will show that in the settings of examples 3.9 and 3.10 definition 3.14 also can be interpreted in graphical terms.

### 3.2.2 Examples: Rating distributions and empirical score distributions

In this section, we look at examples 3.9 and 3.10 in more detail. Observe first that both examples can be described in the same more general terms.

**Assumption G:** There is a finite number of states  $z_1 < z_2 < \dots < z_\ell$  such that

$$\mathbb{P}[S_D \in \{z_1, \dots, z_\ell\}] = 1 = \mathbb{P}[S_N \in \{z_1, \dots, z_\ell\}]. \quad (3.31a)$$

Define for  $i = 1, \dots, \ell$

$$\begin{aligned} \mathbb{P}[S_D = z_i] &= \pi_i, \\ \mathbb{P}[S_N = z_i] &= \omega_i. \end{aligned} \quad (3.31b)$$

**Convention:**

$$z_0 = -\infty. \quad (3.31c)$$

To avoid redundancies in the notation we assume that

$$\pi_i + \omega_i > 0 \text{ for } i \geq 1. \quad (3.31d)$$

Choose  $\ell = n$  and  $z_i = i$  to see that (3.31a) (assumption G) is satisfied in the setting of example 3.9. Then it is obvious how to determine the probabilities  $\pi_i$  and  $\omega_i$ .

In case of example 3.10 choose  $\ell$  to be the number of elements of the set (combined sample)  $\{x_1, \dots, x_{n_D}, y_1, \dots, y_{n_N}\}$  and  $z_i$  as the  $i$ -th element of the ordered list of the different elements of the set. In this case we will have  $1 \leq \ell \leq n_D + n_N$ . The lower extreme case will occur when both the defaulter score sample and the survivor score sample are constant and have the same value. This seems unlikely to happen in practice. The greater limit for  $\ell$  will be assumed when all the values in both the defaulter score and the survivor score samples are pairwise different. This will occur even with probability one if both conditional score distributions are continuous.

For the probabilities  $\pi_i$  and  $\omega_i$  in (3.31b), in the setting of example 3.10 we obtain

$$\begin{aligned} \pi_i &= \delta_{z_i}(x_1, \dots, x_{n_D}) - \delta_{z_{i-1}}(x_1, \dots, x_{n_D}), \\ \omega_i &= \delta_{z_i}(y_1, \dots, y_{n_N}) - \delta_{z_{i-1}}(y_1, \dots, y_{n_N}). \end{aligned} \quad (3.32)$$

**ROC, ROC\*, AUC, and AUC\*.** Under (3.31a) (assumption G), the ROC and ROC\* functions according to (3.5a) and (3.25a) can be described more specifically as follows:

$$\text{ROC}(u) = \begin{cases} 0, & \text{if } 0 = u, \\ \sum_{j=1}^i \pi_j, & \text{if } \sum_{j=1}^{i-1} \omega_j < u \leq \sum_{j=1}^i \omega_j \\ & \text{for } 1 \leq i \leq \ell. \end{cases} \quad (3.33a)$$

$$\text{ROC}^*(u) = \begin{cases} 0, & \text{if } 0 = u, \\ \pi_i/2 + \sum_{j=1}^{i-1} \pi_j, & \text{if } \sum_{j=1}^{i-1} \omega_j < u \leq \sum_{j=1}^i \omega_j \\ & \text{for } 1 \leq i \leq \ell. \end{cases} \quad (3.33b)$$

**Remark 3.17** Observe that equations (3.33a) and (3.33b) can become redundant to some extent in so far as the intervals on their right-hand sides may be empty. This will happen in particular in the context of example 3.10 whenever the samples  $x_1, \dots, x_{n_D}$  and  $y_1, \dots, y_{n_N}$  are disjoint. Let  $\tilde{y}_1 < \dots < \tilde{y}_{k_N}$  be the ordered elements of the set  $\{y_1, \dots, y_{n_N}\}$  of survivor scores. Define  $\tilde{y}_0 = -\infty$ . More efficient versions of (3.33a) and (3.33b) then can be stated as

$$\text{ROC}(u) = \begin{cases} 0, & \text{if } 0 = u, \\ \delta_{\tilde{y}_k}(x_1, \dots, x_{n_D}), & \text{if } \delta_{\tilde{y}_{k-1}}(y_1, \dots, y_{n_N}) < u \leq \delta_{\tilde{y}_k}(y_1, \dots, y_{n_N}) \\ & \text{for } 1 \leq k \leq \ell. \end{cases}$$

$$\text{ROC}^*(u) = \begin{cases} 0, & \text{if } 0 = u, \\ \delta_{\tilde{y}_k}^*(x_1, \dots, x_{n_D}), & \text{if } \delta_{\tilde{y}_{k-1}}(y_1, \dots, y_{n_N}) < u \leq \delta_{\tilde{y}_k}(y_1, \dots, y_{n_N}) \\ & \text{for } 1 \leq k \leq \ell. \end{cases}$$

Under (3.31a) (assumption G) we obtain for the set gROC from definition 3.1

$$\text{gROC} = \left\{ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \omega_1 \\ \pi_1 \end{pmatrix}, \begin{pmatrix} \omega_1 + \omega_2 \\ \pi_1 + \pi_2 \end{pmatrix}, \dots, \begin{pmatrix} \sum_{j=1}^{\ell-1} \omega_j \\ \sum_{j=1}^{\ell-1} \pi_j \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\}. \quad (3.34)$$

Under assumption (3.31d), the points in gROC will be pairwise different. Hence there won't be any redundancy in the representation (3.34) of gROC.

As both the graphs of the ROC and the ROC\* functions as specified by (3.33a) and (3.33b) can obviously be discontinuous at  $u = 0, u = \omega_1, \dots, u = \sum_{j=1}^{\ell-1} \omega_j$ , in practice (see, e.g., Newson, 2001; Fawcett, 2004; Engelmann et al., 2003b) they are often replaced by the linearly interpolated graph through the points of the set gROC as given by (3.34) (in the order of the points as listed there).

**Proposition 3.18** Under (3.31a) (assumption G), the area in the Euclidean plane enclosed by the  $x$ -axis, the vertical line through  $x = 1$  and the graph defined by linear interpolation of the ordered point set gROC as given by (3.34) equals AUC\* as defined by (3.27a) and (3.33b). Moreover, AUC\* can be calculated as

$$\text{AUC}^* = 1/2 \sum_{i=1}^{\ell} \omega_i \pi_i + \sum_{i=2}^{\ell} \omega_i \sum_{j=1}^{i-1} \pi_j. \quad (3.35a)$$

**Proof.** Engelmann et al. (2003b, section III.1.2) showed that the area under the interpolated ROC curve equals AUC\* as represented by (3.28a). Equation (3.35a) follows immediately from (3.28a) and (3.31b).  $\square$

Still under (3.31a) (assumption G), it is easy to see that AUC from definition 3.5, i.e. the “continuous” version of the area under the curve, can be calculated as

$$\text{AUC} = \sum_{i=2}^{\ell} \omega_i \sum_{j=1}^i \pi_j \geq \text{AUC}^*. \quad (3.35b)$$

Observe that  $\text{AUC} = \text{AUC}^*$  if and only if  $\sum_{i=1}^{\ell} \omega_i \pi_i = \text{P}[S_D = S_N] = 0$ .

**Remark 3.19** *In the specific setting of example 3.10, the representation of  $\text{ROC}^*(u)$  from remark 3.17 implies*

$$\text{AUC}^* = (n_D n_N)^{-1} \sum_{i=1}^{n_D} \sum_{j=1}^{n_N} \delta_{y_j}^*(x_i). \quad (3.36a)$$

*The right-hand side of (3.36a) is up to the factor  $n_D n_N$  identical to the statistic of the Mann-Whitney test on whether a distribution is stochastically greater than another distribution (see, e.g., Engelmann et al., 2003b). By means of the representation of  $\text{ROC}(u)$  from remark 3.17, it is not either hard to show that*

$$\text{AUC} = (n_D n_N)^{-1} \sum_{i=1}^{n_D} \sum_{j=1}^{n_N} \delta_{y_j}(x_i). \quad (3.36b)$$

*Clearly,  $\text{AUC}^* = \text{AUC}$  if and only if the samples  $x_1, \dots, x_{n_D}$  and  $y_1, \dots, y_{n_N}$  are disjoint.*

**CAP, CAP\*, AR, and AR\*.** Recall from (2.2) that  $p$  stands for the unconditional probability of default<sup>5</sup>. Under (3.31a) (assumption G), (3.31b) therefore implies that  $\text{P}[S = z_i] = p \pi_i + (1 - p) \omega_i$ . With this in mind, the following representations of  $\text{CAP}(u)$  and  $\text{CAP}^*(u)$  are obvious:

$$\text{CAP}(u) = \begin{cases} 0, & \text{if } 0 = u, \\ \sum_{j=1}^i \pi_j, & \text{if } \sum_{j=1}^{i-1} (p \pi_j + (1 - p) \omega_j) < u \leq \sum_{j=1}^i (p \pi_j + (1 - p) \omega_j) \\ & \text{for } 1 \leq i \leq \ell. \end{cases} \quad (3.37a)$$

$$\text{CAP}^*(u) = \begin{cases} 0, & \text{if } 0 = u, \\ \pi_i/2 + \sum_{j=1}^{i-1} \pi_j, & \text{if } \sum_{j=1}^{i-1} (p \pi_j + (1 - p) \omega_j) < u \leq \sum_{j=1}^i (p \pi_j + (1 - p) \omega_j) \\ & \text{for } 1 \leq i \leq \ell. \end{cases} \quad (3.37b)$$

Note that thanks to assumption (3.31d) the redundancy issue mentioned in remark 3.17 will not occur for representations<sup>6</sup> (3.37a) and (3.37b).

Under (3.31a) (assumption G) we obtain for the set gCAP from definition 3.1

$$\text{gCAP} = \left\{ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} p \pi_1 + (1 - p) \omega_1 \\ \pi_1 \end{pmatrix}, \dots, \begin{pmatrix} \sum_{j=1}^{\ell-1} (p \pi_j + (1 - p) \omega_j) \\ \sum_{j=1}^{\ell-1} \pi_j \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\}. \quad (3.38)$$

As the both the graphs of the CAP and the CAP\* functions as specified by (3.37a) and (3.37b) are obviously discontinuous at  $u = 0, u = p \pi_1 + (1 - p) \omega_1, \dots, u = \sum_{j=1}^{\ell-1} (p \pi_j + (1 - p) \omega_j)$ , in practice (see, e.g., Engelmann et al., 2003b) they are often replaced by the linearly interpolated graph through the points of the set gCAP as given by (3.38) (in the order of the points as listed there).

**Proposition 3.20** *Under (3.31a) (assumption G), the ratio of 1) the area in the Euclidean plane enclosed by the line  $x = y$ , the vertical line through  $x = 1$  and the graph defined by linear*

<sup>5</sup>In example 3.9, the value of  $p$  is a model parameter that can be chosen as it is convenient. In contrast, in example 3.10 a natural (but not necessary) choice for the value of  $p$  is  $p = \frac{n_D}{n_D + n_N}$ .

<sup>6</sup>For more efficient calculations of  $\text{CAP}(u)$  or  $\text{CAP}^*(u)$  in the setting of example 3.10 nonetheless the observation might be useful that  $\sum_{j=1}^i (p \pi_j + (1 - p) \omega_j) = \delta_{z_i}(x_1, \dots, x_{n_D}, y_1, \dots, y_{n_N})$  if  $p$  is chosen as suggested in footnote 5.

interpolation of the ordered point set gCAP as given by (3.38) and 2) the area enclosed by the line  $x = y$ , the vertical line through  $x = 1$  and the CAP\* curve of a perfect score function equals  $\text{AR}^*$  as defined by (3.27b) and (3.37b). Moreover,  $\text{AR}^*$  can be calculated as

$$\text{AR}^* = \sum_{i=1}^{\ell} \omega_i \pi_i + 2 \sum_{i=2}^{\ell} \omega_i \sum_{j=1}^{i-1} \pi_j - 1. \quad (3.39a)$$

**Proof.** As in Engelmann et al. (2003b, section III.1.2) one can show that the area under the interpolated CAP curve equals  $\text{P}[S_D < S] + \text{P}[S_D = S]/2$  where  $S_D$  and  $S$  are independent random variables with the empirical distribution of the scores conditional on default and the unconditional empirical score distribution, respectively. If  $S_N$  denotes a further independent random variable, with the distribution of the scores conditional on survival, and  $S'_D$  is an independent copy of  $S_D$ , this observation implies that

$$\begin{aligned} \text{Ratio of the areas 1) and 2)} &= \frac{\text{P}[S_D < S] + \text{P}[S_D = S]/2 - 1/2}{1 - p/2 - 1/2} \\ &= \frac{2}{1-p} (p \text{P}[S_D < S'_D] + (1-p) \text{P}[S_D < S_N] \\ &\quad + \text{P}[S_D = S'_D]/2 + (1-p) \text{P}[S_D = S_N]/2) \\ &= 2 \text{P}[S_D < S_N] + \text{P}[S_D = S_N] - 1. \end{aligned}$$

By proposition 3.15, this implies the first part of the assertion. (3.39a) then is an immediate consequence of (3.35a) and proposition 3.15 once again.  $\square$

Still under (3.31a) (assumption G), by proposition 3.11 one can conclude that AR from definition 3.5, i.e. the “continuous” version of the accuracy ratio, can be calculated as

$$\text{AR} = 2 \sum_{i=2}^{\ell} \omega_i \sum_{j=1}^i \pi_j - 1 + \frac{p}{1-p} \sum_{i=1}^{\ell} \pi_i^2 > \text{AR}^*. \quad (3.39b)$$

The “ $>$ ” on the right-hand side of (3.39b) is implied by (3.31a) (i.e. at least one  $\pi_i$  is positive).

**Remark 3.21** *In the specific setting of example 3.10, equation (3.39a) is equivalent to*

$$\text{AR}^* = \frac{2}{n_D n_N} \sum_{i=1}^{n_D} \sum_{j=1}^{n_N} \delta_{y_j}^*(x_i) - 1. \quad (3.40a)$$

If  $p = \frac{n_D}{n_D + n_N}$ , by combining proposition 3.11 and (3.36b) one can also calculate AR for the setting of example 3.10, i.e. a representation equivalent to (3.39b):

$$\text{AR} = \frac{2}{n_D n_N} \sum_{i=1}^{n_D} \sum_{j=1}^{n_N} \delta_{y_j}(x_i) - 1 + \frac{1}{n_D n_N} \sum_{i=1}^{n_D} \sum_{j=1}^{n_D} \delta_{x_j}(x_i) \delta_{x_i}(x_j). \quad (3.40b)$$

Note that  $\text{AR} > \text{AR}^*$  even if the samples  $x_1, \dots, x_{n_D}$  and  $y_1, \dots, y_{n_N}$  are disjoint. This follows from  $\sum_{i=1}^{n_D} \sum_{j=1}^{n_D} \delta_{x_j}(x_i) \delta_{x_i}(x_j) \geq \sum_{i=1}^{n_D} 1 = n_D > 0$ .

## 4 Discriminatory power: Numerical aspects

Engelmann et al. (2003a,b) compared for different sample sizes approximate normality-based and bootstrap confidence intervals for AUC. As they worked with a huge dataset of defaulter and non-defaulter scores, they treated the estimates on the whole dataset as “true” values – an assumption confirmed by tight confidence intervals. Engelmann et al. then sub-sampled from the dataset to study the impact of smaller sample sizes. Their conclusion – for scores both on continuous and discrete scales – was that even for defaulter samples of size ten the approximate and bootstrap intervals do not differ much and cover the “true” value.

After having presented some general considerations on the impact on bootstrap performance by sample size in sub-section 4.1, in sections 4.2 and 4.3 we supplement the observations of Engelmann et al. in a simulation study<sup>7</sup> where we sample from known analytical distributions. This way, we really know the true value of AUC and can determine whether or not the true value is covered by a confidence interval. Additionally, we study the impact of having an even smaller sample size of five defaulters.

Note that by proposition 3.15 any conclusion on estimation uncertainty for AUC\* also applies to AR\*.

### 4.1 Bootstrap confidence intervals when the default sample size is small

Davison and Hinkley (1997, section 2.3) commented on the question of how large the sample size should be in order to generate meaningful bootstrap samples. Davison and Hinkley observed that if the size of the original sample is  $n$  the number of different bootstrap samples that can be generated from this sample is no larger than  $\binom{2n-1}{n}$ . Table 1 shows the value of this term for the first eleven positive integers. When following the general recommendation by Davison and Hinkley to generate at least 1000 bootstrap samples, according to table 1 then beginning with  $n = 7$  it is possible not to have any identical (up to permutations) samples. For sample size six and below the sample variation will be restricted for combinatorial reasons. This applies even more to samples on a discrete scale which in most cases include ties. One should therefore expect that bootstrap intervals for AUC become less reliable when the size of the defaulter score sample is six or less or when the sample includes ties. A simple simulation experiment further illustrates this observation. For two samples of size  $n \in \{1, \dots, 11\}$  with  $n$  different elements and  $n - 1$  different elements respectively, we run<sup>8</sup> 100 bootstrap experiments each with 1000 iterations. In each bootstrap experiment we count how many of the generated samples are different.

Table 2 indeed clearly demonstrates that the factual sample size from bootstrapping is significantly smaller than the nominal bootstrap sample size when the original sample has less than nine elements. The impact of small size of the original sample is even stronger when the original sample includes at least one tie (two identical elements). Observe, however, that the impact of diminished factual sample size is partially mitigated by the fact that for combinatorial reasons the frequencies of duplicated bootstrap samples will have some variation.

---

<sup>7</sup>Like Engelmann et al. (2003a,b) we compare approximate normality-based and bootstrap confidence intervals for AUC. Newson (2006) describes how jackknife methods can be applied to estimate confidence intervals for Somers’ D (and hence in particular for AUC).

<sup>8</sup>All calculations for this paper were conducted with R version 2.6.2 (R Development Core Team, 2008).

**Bootstrap confidence intervals.** In sections 4.2 and 4.3 we calculate *basic bootstrap intervals* generated by *nonparametric bootstrap* as described in section 2.4 of Davison and Hinkley (1997). Technically speaking, if the original estimate of a parameter (e.g. of  $AUC^*$ ) is  $t$  and we have a bootstrap sample  $t_1^* \leq t_2^* \leq \dots \leq t_n^*$  of estimates for the same parameter, then the basic bootstrap interval  $I$  at confidence level  $\gamma \in (0, 1)$  is given by

$$I = [2t - t_{n(1+\gamma)/2}^*, 2t - t_{n(1-\gamma)/2}^*], \quad (4.1)$$

where we assume that  $(n+1)(1+\gamma)/2$  and  $(n+1)(1-\gamma)/2$  are integers in the range from 1 to  $n$ . Our standard choice of  $n$  and  $\gamma$  in sections 4.2 and 4.3 is  $n = 999$  and  $\gamma = 95\%$ , leading to  $(n+1)(1+\gamma)/2 = 975$  and  $(n+1)(1-\gamma)/2 = 25$ .

**Approximate confidence intervals for  $AUC^*$  based on the central limit theorem.** Additionally, in sections 4.2 and 4.3 we calculate *approximate confidence intervals* for  $AUC$  according to Engelmann et al. (2003b, equation (12)).

## 4.2 Simulation study: Continuous score distributions

We consider the normal distribution example from section 3.1.1 with the following choice of parameters:

$$\begin{aligned} \mu_D &= 6.8, \sigma_D = 1.96 \\ \mu_N &= 8.5, \sigma_N = 2 \end{aligned} \quad (4.2)$$

These parameters are chosen such as to match the first two moments of the binomial distributions looked at in subsequent section 4.3. According to (3.14), under the normal assumption with parameters as in (4.2) then we have  $AUC = AUC^* = 71.615\%$ .

For defaulter score sample sizes  $n_D \in \{5, 10, 15, 20, 25, 30, 35, 40, 45, 50\}$  and constant survivor score sample size  $n_N = 250$ , we conduct  $k = 100$  times the following bootstrap experiment:

- 1) Simulate a sample of size  $n_D$  of independent normally distributed defaulter scores and a sample of size  $n_N$  of independent normally distributed survivor scores, with parameters as specified in (4.2).
- 2) Based on the samples from step 1) calculate estimates  $AUC_{\text{kernel}}$  according to (3.18) and (3.19c) and  $AUC_{\text{emp}}$  according to (3.36a) for  $AUC$ .
- 3) Based on the samples from step 1) and  $AUC_{\text{emp}}$  calculate the normal 95% confidence interval  $I_{\text{normal}}$  (as described by Engelmann et al., 2003b, equation (12)).
- 4) Generate for each of the two samples from step 1)  $r = 999$  nonparametric bootstrap samples, thus obtaining  $r = 999$  pairs of bootstrap samples.
- 5) For each pair of bootstrap samples associated with bootstrap trial  $i = 1, \dots, r$  calculate estimates  $\widehat{AUC}_i$  according to (3.18) and (3.19c) as well as  $\widetilde{AUC}_i$  according to (3.36a) for  $AUC$ .
- 6) Calculate basic bootstrap 95% confidence intervals  $I_{\text{kernel}}$  and  $I_{\text{emp}}$  as described in (4.1) based on the estimate  $AUC_{\text{kernel}}$  and the sample  $(\widehat{AUC}_i)_{i=1, \dots, r}$  and the estimate  $AUC_{\text{emp}}$  and the sample  $(\widetilde{AUC}_i)_{i=1, \dots, r}$ , respectively.

7) Check whether or not

$$\begin{array}{lll} \text{AUC} \in I_{\text{normal}}, & \text{AUC} \in I_{\text{kernel}}, & \text{AUC} \in I_{\text{emp}}, \\ 50\% \in I_{\text{normal}}, & 50\% \in I_{\text{kernel}}, & 50\% \in I_{\text{emp}}. \end{array}$$

To give an impression of the variation encountered with the different confidence interval methodologies and the different sample sizes, table 3 (for defaulter sample sizes  $n_D = 5$ ,  $n_D = 25$ , and  $n_D = 45$ ) shows the AUC estimates from the original samples and the related confidence interval estimates for the first five experiments. Although it is clear from the tables that the estimates are more stable and the confidence intervals are tighter for the larger defaulter score samples, it is nonetheless hard to conclude from these results which of the estimation methods is most efficient.

Table 4 and figure 6 therefore provide information on how often the true AUC was covered by the confidence intervals and how often 50% was an element of the confidence intervals. The check of the coverage of 50% is of interest because as long as 50% is included in a 95% confidence interval for AUC, one cannot conclude that the score function or rating system under consideration has got any discriminatory power.

According to table 4 and figure 6 coverage of the true AUC is poor for defaulter sample size  $n_D \leq 15$  but becomes satisfactory for the larger defaulter sample sizes. At the same time, the values of coverage of 50% indicate poor power for defaulter sample sizes  $n_D \leq 20$  and much better power for defaulter sample size  $n_D = 25$  and larger.

For all defaulter sample sizes the coverage differences both for true AUC and for 50% are negligible in case of the “empirical” confidence intervals and the kernel estimation-based confidence intervals. For the smaller defaulter sample sizes ( $n_D \leq 15$ ), coverage of true AUC by the normal confidence interval is clearly better than by the “empirical” confidence intervals and the kernel estimation-based confidence intervals but still less than the nominal level of 95%. The better coverage of true AUC by the normal confidence intervals, however, comes at the price of a much higher coverage of 50% for defaulter samples sizes  $n_D \leq 20$  (type II error). For defaulter sample sizes  $n_D \geq 25$  differences in performance of the three approaches to confidence intervals seem to vanish.

**Remark 4.1** *With a view on (3.36a), it follows from the duality of tests and confidence intervals (see, e.g., Casella and Berger, 2002, theorem 9.2.2) that the check of whether 50% is covered by the AUC 95% confidence interval is equivalent to conducting a Mann-Whitney test of whether the defaulter score distribution and the survivor score distribution are equal (null hypothesis). The exact distribution of the Mann-Whitney test statistic can be calculated with standard statistical software packages. Hence the 95% confidence interval coverage rates of 50% reported in table 4 can be double-checked against type II error rates from application of the two-sided Mann-Whitney test at 5% type I error level.*

The type II error rates mentioned in remark 4.1 are displayed in the second to last column of table 4. For the sake of completeness, in the last column of table 4 type II error rates from application of the two-sided Kolmogorov-Smirnov test are presented, too. Comparison of the Mann-Whitney type II error and the coverage of 50% by the AUC confidence intervals clearly indicates that for defaulter sample size  $n_D \leq 20$  the bootstrap confidence intervals are too narrow. With a view on table 2 this observation does not come as a surprise for very small



defaulter sample sizes but is slightly astonishing for a defaulter sample size like  $n_D = 15$ . The confidence intervals based on asymptotic normality, however, seem to perform quite well for sample size  $n_D \geq 10$ . Comparing the last column of table 4 to the second-last column moreover shows that the Mann-Whitney test is clearly more powerful for smaller defaulter sample sizes than the Kolmogorov-Smirnov test.

In summary, the simulation results suggest that in the continuous setting of this section for defaulter sample size  $n_D \geq 20$  the performance differences between the three approaches to AUC confidence intervals considered are negligible. For defaulter sample size  $n_D < 20$ , however, with a view on the coverage of the true AUC parameter it seems clearly preferable to deploy the confidence interval approach based on asymptotic normality (as described, e.g., by Engelmann et al., 2003a,b) because its coverage rates come closest to the nominal confidence level (but are still smaller). For very small defaulter sample size  $n_D \leq 10$ , poorer coverage of the true AUC parameter may come together with a high type II error (high coverage of 50%, indicating misleadingly that the score function is powerless).

On the basis of a more intensive simulation study that includes observations on coverage rates, thus we can re-affirm and at the same time refine the conclusion by Engelmann et al. (2003a,b) that confidence intervals for AUC (and AR) based on asymptotic normality work reasonably well for sample data on a continuous scale, even for small defaulter sample size like  $n_D = 10$  but not necessarily for a very small defaulter sample size like  $n_D = 5$ . Moreover, for defaulter sample sizes  $n_D < 20$  the asymptotic normality confidence interval estimator out-performs bootstrap-based estimators.

### 4.3 Simulation study: Discrete score distributions

We consider the binomial distribution example for 17 rating grades from figure 3 with probability parameter  $p_D = 0.4$  for the defaulter rating distribution and probability parameter  $p_N = 0.5$  for the survivor rating distribution. As a consequence, the first two moments of the defaulter rating distribution match the first two moments of the defaulter score distribution from section 4.2 and the first two moments of the survivor rating distribution match the first two moments of the survivor score distribution from section 4.2. Moreover, also the discriminatory power of the fictitious rating system considered in this section is almost equal to the discriminatory power of the score function from section 4.2 (AUC\* 71.413% according to (3.35a) vs. AUC 71.615%).

To assess the impact of the discreteness of the model, we conduct the same simulation exercise as in section 4.2 but replace step 1) by step 1\*) which reads

- 1\*) Simulate a sample of size  $n_D$  of independent binomially distributed defaulter ratings and a sample of size  $n_N$  of independent binomially distributed survivor ratings, with probability parameters  $p_D = 0.4$  and  $p_N = 0.5$  respectively.

As in section 4.2, to give an impression of the variation encountered with the different confidence interval methodologies and the different sample sizes, table 5 (for defaulter sample sizes  $n_D = 5$ ,  $n_D = 25$ , and  $n_D = 45$ ) shows the AUC estimates from the original samples and the related confidence interval estimates for the first five experiments. Although it is clear from the tables that the estimates are more stable and the confidence intervals are tighter for the larger defaulter score samples, it is nonetheless hard to conclude from these results which of the estimation methods is most efficient. Interesting is also the result from experiment number 5 for sample

size  $n_D = 5$  in table 5 which with lower confidence bounds of 90.0% and more looks very much like an outlier due to a defaulter sample concentrated at the bad end of the rating scale.

Table 6 and figure 7 provide information on how often the true AUC was covered by the confidence intervals and how often 50% was an element of the confidence intervals. In contrast to table 4 and figure 6, table 6 and figure 7 do not give a very clear picture of the performance of the three AUC estimation approaches on the rating data. While coverage of 50% (type II error) is high for defaulter sample sizes smaller than  $n_D = 30$ , coverage of 50% reaches very small values as in the continuous case of section 4.2 for larger defaulter sample sizes. Presumably due to the relatively small number of 100 bootstrap experiments – which already requires some hours of computation time –, according to figure 7 there is some variation and not really a clear trend in the level of coverage of the true AUC parameter. Even for a relatively high defaulter sample size of  $n_D = 40$  there is sort of a collapse of coverage of true AUC with percentages of 90% or lower. For defaulter sample size of  $n_D = 45$  or more there might be some stabilisation at a satisfactory level.

As in the continuous case, for all defaulter sample sizes the coverage differences both for true AUC and for 50% are negligible in case of the “empirical” confidence intervals and the kernel estimation-based confidence intervals. For the smaller defaulter sample sizes ( $n_D \leq 20$ ), coverage of true AUC by the normal confidence interval is clearly better than by the “empirical” confidence intervals and the kernel estimation-based confidence intervals but still less than the nominal level of 95%. The better coverage of true AUC by the normal confidence intervals, however, comes at the price of a much higher coverage of 50% for defaulter samples sizes  $n_D \leq 15$  (type II error). For defaulter sample sizes  $n_D \geq 25$  differences in performance of the three approaches to confidence intervals seem to vanish.

**Remark 4.2** *Remark 4.1 essentially also applies to the setting of this section. But take into account that in the presence of ties in the sample the equivalence between  $AUC^*$  as defined by (3.27a) and the Mann-Whitney statistic only holds when ranks for equal elements of the ordered total sample are assigned as mid-ranks. With this in mind we can double-check the 95% confidence coverage rates of 50% reported in table 6 against type II error rates from application of the two-sided Mann-Whitney test at 5% type I error level in the same manner as we have done for remark 4.1.*

The type II error rates<sup>9</sup> mentioned in remark 4.2 are reported in the second to last column of table 6. We have presented type II error rates from application of the two-sided Kolmogorov-Smirnov test in the last column of table 4. Due to the massive presence of ties in the discrete-case samples, however, application of the Kolmogorov-Smirnov test does not seem appropriate in this section. Instead, we report type II error rates from application of the two-sided exact Fisher test<sup>10</sup> (see, e.g., Weisstein, 2009) in the last column of table 6.

Again, comparison of the Mann-Whitney type II error and the coverage of 50% by the AUC confidence intervals clearly indicates that for defaulter sample size  $n_D \leq 20$  the bootstrap confidence intervals are too narrow. With another view on table 2 this is even less a surprise than in the continuous case. The confidence intervals based on asymptotic normality, however,

---

<sup>9</sup>Exact p-values for the Mann-Whitney test on samples with ties were calculated with the function `wilcox.test` from the R-software package `coin`.

<sup>10</sup>The p-values of Fisher’s exact test have been calculated with the function `fisher.test` (R-software package `stats`) in simulation mode due to too high memory and time requirements of the exact mode.

seem again to perform quite well for sample size  $n_D \geq 10$ . Comparing the last column of table 4 to the second-last column moreover shows that the Mann-Whitney test is clearly more powerful for smaller and even some moderate defaulter sample sizes than Fisher’s test.

In summary, the simulation results suggest that in the discrete setting of this section for defaulter sample size  $n_D \geq 25$  the performance differences between the three approaches to AUC confidence intervals considered are negligible. For defaulter sample size  $n_D \leq 20$ , however, with a view on the coverage of the true AUC parameter it seems clearly preferable to deploy the confidence interval approach based on asymptotic normality (as described, e.g., by Engelmann et al., 2003a,b) because its coverage rates come closest to the nominal confidence level (but are still smaller). For very small defaulter sample size  $n_D \leq 10$ , poorer coverage of the true AUC parameter may come together with a high type II error (high coverage of 50%, indicating misleadingly that the score function is powerless).

On the basis of this more intensive simulation study that includes observations on coverage rates, thus we can re-affirm and at the same time refine the interesting conclusion by Engelmann et al. (2003a,b) that confidence intervals for AUC (and AR) based on asymptotic normality work reasonably (when compared to other approaches) for sample data on a discrete scale, even for small defaulter sample size like  $n_D = 10$  (but not necessarily for smaller defaulter sample sizes). Bootstrap estimators are out-performed for such small sample sizes by the asymptotic normality estimator. While the performance of the three AUC confidence interval methods does not seem to be much worse for discrete scale rating distributions than for continuous scale score distributions, there might be a higher likelihood of performance outliers – as discussed at the beginning of section 4.1.

## 5 Determining PD curves by parametric estimation of ROC and CAP curves

In the context of portfolios with little default data, van der Burgt (2008) suggested estimating the conditional probability of default by fitting a one-parameter curve to the observed CAP curve, taking the derivative of the fitted curve and then calculating the conditional probabilities of default based on the derivative. Van der Burgt did not provide much background information on why he chose the specific one-parameter family of curves he used in the paper nor did he look more closely at the properties of this family of curves.

In section 5.1, we provide some background information on the implicit assumptions and implications when working with parametric approaches for CAP and ROC functions. In section 5.2, we discuss van der Burgt’s approach in detail and introduce three potential alternatives. In section 5.3 we compare the performance of the four approaches by looking at some numerical examples.

Parametric approaches (e.g. logit or van der Burgt’s approaches) to PD curves are popular with practitioners because they can be designed such as to guarantee monotonicity of the conditional PD estimates. With an appropriate choice of the parametric shape, it is also possible to replicate the exponential-like growth that is observed for corporate default rates associated with agency ratings. However, as discussed in Tasche (2008, section 4), monotonicity of PD curves must not be taken for granted. We will see in section 5.3 that both the erroneous assumption of monotonicity and the mis-specification of discriminatory power can cause huge estimation errors

when following a parametric approach to PD curve estimation. See Pluto and Tasche (2005) for a non-parametric approach that might be a viable alternative in particular when little default observation is available.

## 5.1 Derivatives of CAP and ROC curves

It is a well known fact that there is a close link between ROC and CAP curves on the one hand and conditional probabilities of default on the other hand. Technically speaking, the link is based on the following easy-to-prove (when making use of theorem 3.2) observation.

**Proposition 5.1** *Let  $F_D$  and  $F_N$  be distribution functions on an open interval  $I \subset \mathbb{R}$ . Assume that  $F_D$  has a density  $f_D$  which is continuous on  $I$  and that  $F_N$  has a positive density  $f_N$  that is continuous on  $I$ . Let  $0 < p < 1$  be a fixed probability and define the mixed distribution  $F$  by (2.4). Write  $f$  for the density of  $F$ . Define  $\text{ROC}(u)$  and  $\text{CAP}(u)$ ,  $u \in (0, 1)$  by (3.5a) and (3.5b), respectively. Then both ROC and CAP are continuously differentiable for  $u \in (0, 1)$  with derivatives*

$$\text{ROC}'(u) = \frac{f_D(F_N^{-1}(u))}{f_N(F_N^{-1}(u))}, \quad (5.1a)$$

$$\begin{aligned} \text{CAP}'(u) &= \frac{f_D(F^{-1}(u))}{p f_D(F^{-1}(u)) + (1-p) f_N(F^{-1}(u))} \\ &= \frac{f_D(F^{-1}(u))}{f(F^{-1}(u))}. \end{aligned} \quad (5.1b)$$

Proposition 5.1 is of high interest in the context of individual default risk analysis because – in the notation of sections 2.2 and 3 – the probability of default conditional on a score value  $s$  is given by (2.6b). Proposition 5.1 then immediately implies

$$\text{P}[D | S = s] = \frac{p \text{ROC}'(F_N(s))}{p \text{ROC}'(F_N(s)) + 1 - p} \quad (5.2a)$$

$$\begin{aligned} &= p \text{CAP}'(p F_D(s) + (1-p) F_N(s)) \\ &= p \text{CAP}'(F(s)). \end{aligned} \quad (5.2b)$$

Note that by (5.2b), the derivative of a differentiable CAP curve for a borrower population with unconditional probability of default  $p > 0$  is necessarily bounded from above by  $1/p$ . The following theorem shows on the one hand that this condition is not only a necessary but also a sufficient condition for a distribution function on the unit interval to be a CAP curve. On the other hand, the theorem shows that a CAP curve relates not only to one combination of conditional and unconditional score distributions but provides a link between conditional and unconditional score distributions which applies to an infinite number of such combinations.

**Theorem 5.2** *Let  $p \in (0, 1)$  be a fixed probability. Let  $f_D \geq 0$  be a density on  $\mathbb{R}$  such that the set  $I = \{f_D > 0\}$  is an open interval and  $f_D$  is continuous in  $I$ . Denote by  $F_D(s) = \int_{-\infty}^s f_D(v) dv$  the distribution function associated with  $f_D$ . Then the following two statements are equivalent:*

- (i) *The function  $u \mapsto C(u)$ ,  $u \in (0, 1)$  is continuously differentiable in  $u$  with  $\lim_{u \rightarrow 0} C(u) = 0$ ,  $\lim_{u \rightarrow 1} C(u) = 1$ , and  $0 < C'(u) \leq 1/p$ ,  $u \in (0, 1)$ .*

(ii) There is a density  $f_N \geq 0$  such that  $\{f_N > 0\} = I$ ,  $f_N$  is continuous in  $I$  and

$$C(u) = F_D(F^{-1}(u)), \quad u \in (0, 1), \quad (5.3)$$

where  $F(s) = pF_D(s) + (1-p) \int_{-\infty}^s f_N(v) dv$ .

**Proof.**

(i)  $\Rightarrow$  (ii): By assumption,  $C$  maps  $(0, 1)$  onto  $(0, 1)$  and the inverse  $C^{-1}$  of  $C$  exists. Define  $F(s) = C^{-1}(F_D(s))$ . Then  $F$  is a distribution function with  $\lim_{s \rightarrow \inf I} F(s) = 0$ ,  $\lim_{s \rightarrow \sup I} F(s) = 1$  and density

$$f(s) = F'(s) = \frac{f_D(s)}{C'(F(s))}, \quad s \in I. \quad (5.4)$$

Observe that  $f(s)$  is positive and continuous in  $I$ . Hence the inverse  $F^{-1}$  of  $F$  exists. Let

$$\begin{aligned} F_N(s) &= \frac{F(s) - pF_D(s)}{1-p} \quad s \in \mathbb{R}, \quad \text{and} \\ f_N(s) &= f_D(s) \frac{1/C'(F(s)) - p}{1-p}, \quad s \in I. \end{aligned} \quad (5.5)$$

By (5.4), then  $f_N$  is the continuous derivative of  $F_N$  and is positive in  $I$  by assumption on  $C'$  and  $f_D$ . This implies that  $F_N$  is a distribution function with  $\lim_{s \rightarrow \inf I} F_N(s) = 0$ ,  $\lim_{s \rightarrow \sup I} F_N(s) = 1$  and density  $f_N$ . By construction of  $F$  and  $F_N$ , the functions  $C$ ,  $F_D$ , and  $F$  satisfy (5.3).

(ii)  $\Rightarrow$  (i): By construction,  $F_D(s)$  and  $F(s)$  are distribution functions which converge to 0 for  $s \rightarrow \inf I$  and to 1 for  $s \rightarrow \sup I$ . This implies the limit statements for  $C$ . Equation (5.3) implies that  $C$  is continuously differentiable with derivative

$$0 < C'(u) = \frac{f_D(F^{-1}(u))}{p f_D(F^{-1}(u)) + (1-p) f_N(F^{-1}(u))} \leq 1/p.$$

□

For the sake of completeness, we provide without proof the result corresponding to theorem 5.2 for ROC curves. In contrast to the case of CAP curves, essentially every continuously differentiable and strictly increasing distribution function on the unit interval is the ROC curve for an infinite number of combinations of score distributions conditional of default and survival respectively.

**Proposition 5.3** Let  $f_D \geq 0$  be a density on  $\mathbb{R}$  such that the set  $I = \{f_D > 0\}$  is an open interval and  $f_D$  is continuous in  $I$ . Denote by  $F_D(s) = \int_{-\infty}^s f_D(v) dv$  the distribution function associated with  $f_D$ . Then the following two statements are equivalent:

(i) The function  $u \mapsto R(u)$ ,  $u \in (0, 1)$  is continuously differentiable in  $u$  with  $\lim_{u \rightarrow 0} R(u) = 0$ ,  $\lim_{u \rightarrow 1} R(u) = 1$ , and  $0 < R'(u)$ ,  $u \in (0, 1)$ .

(ii) There is a density  $f_N \geq 0$  such that  $\{f_N > 0\} = I$ ,  $f_N$  is continuous in  $I$  and

$$R(u) = F_D(F_N^{-1}(u)), \quad u \in (0, 1), \quad (5.6)$$

where  $F_N(s) = \int_{-\infty}^s f_N(v) dv$ .

The basic idea both with theorem 5.2 and proposition 5.3 is that if in the functional equation  $f(x) = g(h^{-1}(x))$  two of the three functions  $f$ ,  $g$  and  $h$  are given then the third can be calculated by solving the equation for it. In the cases of ROC and CAP curves, matters can get more complicated because the involved functions are not necessarily invertible. This would entail some technicalities when trying to solve  $f(x) = g(h^{-1}(x))$  for  $g$  or  $h$ . However, to relate conditional probabilities of default to ROC and CAP functions via (5.2a) and (5.2b) we need the existence of densities. This introduces some degree of smoothness as can be seen from theorem 5.2 and proposition 5.3. Both the theorem and the proposition could also be stated with fixed distribution  $F_N$  of the survivor scores. However, the survivor score distribution appears in the CAP function only as a mixture with the defaulter score distribution. Therefore, stating theorem 5.2 with given survivor score distribution would no longer be straight-forward and the proof would involve the implicit function theorem. As the additional insight by such a version of theorem 5.2 would be limited, in this paper the formulation of the theorem as provided above has been preferred.

## 5.2 Van der Burgt's approach and alternatives

The one-parameter curve proposed by van der Burgt (2008) for estimating CAP functions is

$$C_\kappa(u) = \frac{1 - e^{-\kappa u}}{1 - e^{-\kappa}}, \quad u \in [0, 1], \quad (5.7a)$$

where  $\kappa \in \mathbb{R}$  is the fitting parameter. The function  $C_\kappa$  is obviously a distribution function on  $[0, 1]$ . Moreover, for positive  $\kappa$  the graph of  $C_\kappa$  is concave as one might expect from the CAP curve of a score function that assigns low scores to bad borrowers and high scores to good borrowers. For  $\kappa \rightarrow 0$  the graph of  $C_\kappa$  converges toward the diagonal line, i.e. the graph of a powerless score function. The derivative of  $C_\kappa$  and  $\text{AR}_\kappa$  associated with  $C_\kappa$  according to (3.10b) are easily computed as

$$C'_\kappa(u) = \frac{\kappa e^{-\kappa u}}{1 - e^{-\kappa}}, \quad (5.7b)$$

$$\text{AR}_\kappa = \frac{2}{1-p} \left( \frac{1}{1 - e^{-\kappa}} - \frac{1}{\kappa} - 1/2 \right). \quad (5.7c)$$

In (5.7c) the parameter  $p > 0$  denotes the unconditional probability of default of the estimation sample in the sense of section 2.1. Observe from (5.7b) that for  $\kappa > 0$

$$C'_\kappa(1) = \frac{\kappa}{1 - e^{-\kappa}} e^{-\kappa} \leq C'_\kappa(u) \leq \frac{\kappa}{1 - e^{-\kappa}} = C'_\kappa(0), \quad u \in [0, 1]. \quad (5.8a)$$

Theorem 5.2 hence implies

$$\kappa < \frac{\kappa}{1 - e^{-\kappa}} \leq \frac{1}{p}. \quad (5.8b)$$

Given a CAP curve  $\text{CAP}(u)$  to be approximated, in the setting of a continuous score function a natural approach to finding the best fit  $\hat{\kappa}$  would be a least squares procedure as the following

$$\begin{aligned} \hat{\kappa} &= \arg \min_{\kappa > 0} \int_0^1 (\text{CAP}(u) - C_\kappa(u))^2 du \\ &= \arg \min_{\kappa > 0} \mathbb{E} \left[ (F_D(S) - C_\kappa(F(S)))^2 \right]. \end{aligned} \quad (5.9a)$$

In practice, the integration in (5.9a) would have to be replaced by a sample mean. Alternatively, van der Burgt (2008) suggested inferring  $\kappa$  by means of (5.7c) from an estimated<sup>11</sup> AR (or via proposition 3.6 from an estimated AUC). Assuming that an estimate of the unconditional probability  $p$  is available, probabilities of default conditional on realised score values then can be estimated via (5.2b):

$$P[D | S = s] = \frac{p \kappa e^{-\kappa F(s)}}{1 - e^{-\kappa}}. \quad (5.9b)$$

Van der Burgt (2008), however, applied the methodology to a discrete setting as described in example 3.9 (rating system with  $n$  grades). In the notation of example 3.9 van der Burgt's approach to finding the best fit parameter  $\hat{\kappa}$  can be described as

$$\begin{aligned} \hat{\kappa} &= \arg \min_{\kappa > 0} \sum_{j=1}^n (P[R_D \leq j] - C_{\kappa}(P[R \leq j]))^2 \\ &= \arg \min_{\kappa > 0} \mathbb{E} \left[ \frac{(F_D(R) - C_{\kappa}(F(R)))^2}{P[R = r] |_{r=R}} \right]. \end{aligned} \quad (5.10a)$$

Van der Burgt's approach hence can be regarded as sort of an unweighted regression in which the same weights are given to rating grades with very few observations and grades with quite a lot of observations. For calculating the estimates of the conditional probabilities of default van der Burgt (2008) does not deploy equation (5.9b) but a modification that substitutes the unconditional score distribution function  $F$  by the mean of its right and left continuous versions  $(F + F(\cdot - 0))/2$ :

$$P[D | R = j] = \frac{p \kappa \exp(-\kappa (P[R < j] + P[R \leq j])/2)}{1 - e^{-\kappa}}. \quad (5.10b)$$

In his paper, van der Burgt (2008) does not spend much time with explaining the why and how of his approach. It is tempting to guess that the approach was more driven by the results than by theoretical considerations. We observe that it is not obvious how to deploy van der Burgt's regression approach (5.10a) for a sample of scores from a score function with continuous scale. Therefore, for our example calculations in section 5.3 we will make use of (5.9a) and (5.9b) for the continuous setting and of (5.10a) and (5.10b) for the discrete setting of example 3.9.

In general, when choosing  $C_{\kappa}$  for fitting a CAP curve, one should be aware that as a consequence of theorem 5.2 this choice implies some structural links between the score distribution of the defaulters and the score distribution of the survivors. This is illustrated in figure 8 which shows for unconditional probability of default  $p = 0.01$  and different values of  $\kappa$  the survivor score densities that are implied by theorem 5.2 when the defaulter score density is assumed to be standard normal. Clearly, for large  $\kappa$  and, by (5.7c), high discriminatory power the implied survivor score distributions are not normal as they are not symmetric.

This observation on the one hand might be considered not very appealing. On the other hand, it suggests an alternative approach along the lines of section 3.1.1 which provides in (3.14) a two-parametric representation of the ROC function for the case of normally distributed defaulter and survivor score distributions.

As mentioned in section 3.1.1, no closed form is available for the CAP function in case of normally distributed defaulter and survivor score distributions. This is one reason why we consider in the

---

<sup>11</sup>This requires that there is an estimate of  $p$ . Van der Burgt assumes  $p \approx 0$  for the purpose of estimating  $\kappa$  but then, in a further step, makes use of the fact that  $p$  is positive.

following how to approximate general ROC curves (not CAP curves) by the ROC function of the normal family as described in section 3.1.1. Another reason is that, in general, fitting ROC curves is conceptually sounder than fitting CAP curves because this way one can better separate the estimation of an unconditional probability of default from the estimation of parameters of the fitting function.

By substituting in (3.14) the parameter  $b > 0$  for  $\sigma_N/\sigma_D$  and the parameter  $a \in \mathbb{R}$  for  $\frac{\mu_N - \mu_D}{\sigma_D}$ , we obtain a two-parametric family of ROC functions:

$$R_{a,b}(u) = \Phi(a + b\Phi^{-1}(u)), \quad u \in (0, 1). \quad (5.11a)$$

From this, it follows readily that

$$R'_{a,b}(u) = b \frac{\varphi(a + b\Phi^{-1}(u))}{\varphi(\Phi^{-1}(u))}, \quad (5.11b)$$

$$\text{AR}_{a,b} = 2\Phi\left(\frac{a}{\sqrt{b^2 + 1}}\right) - 1. \quad (5.11c)$$

Clearly, a two-parameter family of functions will give better fits than a one-parameter family of functions as the one-parameter family proposed by van der Burgt (2008). In order to have a fair comparison, therefore, in the following we will focus on the one-parameter sub-family of (5.11a) specified by fixing  $b$  at  $b = 1$ . We simplify notation by writing  $R_a$  for  $R_{a,1}$ . From section 3.1.1 it follows that the one-parameter family of ROC functions  $R_a$  includes, in particular, the ROC curves for normally distributed defaulter and survivor score functions when their variances are equal. Equations (5.11b) and (5.11c) are simplified significantly for  $R_a$ :

$$R'_a(u) = e^{-a\Phi^{-1}(u) - a^2/2}, \quad (5.12a)$$

$$\text{AR}_a = 2\Phi(a/\sqrt{2}) - 1. \quad (5.12b)$$

When the unconditional probability of default  $p$  is known, (5.12a) via (5.2a) implies the following representation of the probability of default conditional on a realised score value:

$$\text{P}[D | S = s] = \frac{1}{1 + \frac{1-p}{p} \exp(a\Phi^{-1}(F_N(s)) + a^2/2)}. \quad (5.13a)$$

Clearly, (5.13a) can be rewritten as

$$\text{P}[D | S = s] = \frac{1}{1 + \exp(\alpha + \beta\Phi^{-1}(F_N(s)))} \quad (5.13b)$$

with

$$\alpha = \log\left(\frac{1-p}{p}\right) + a^2/2, \quad \beta = a.$$

Thus, the conditional PDs derived from the one-parameter ROC approximation approach (5.11a) with  $b = 1$  are the conditional PDs of a logit regression where the default indicator is regressed on the explanatory variable  $\Phi^{-1}(F_N(S))$ . In the special case where the score distribution conditional on survival is normal (i.e.  $F_N(s) = \Phi((s - \mu)/\sigma)$  for some suitable constants  $\mu$  and  $\sigma$ ), the right-hand side of equation (5.13b) coincides with the conditional PDs of the common logit approach:

$$\text{P}[D | S = s] = \frac{1}{1 + \exp(\alpha + \beta s)}. \quad (5.14)$$



Thus (5.13b) can be considered a *robust logit approach* that replaces regression on the original score  $S$  by regression on the transformed score  $\Phi^{-1}(F_N(S))$  to account for the fact that the score distribution might not be normal. As such, the suggestion by van der Burgt (2008) leads to a potentially quite useful modification of logit regression in the univariate case.

On an estimation sample in the sense of section 2.1.1, parameters for logit-type raw conditional PDs as specified in equations (5.13b) and (5.14) can be estimated by maximum likelihood (MLE) procedures (see, e.g., Cramer, 2003, chapter 3). In the case of (5.13b), MLE will only work if  $0 < F_N(x_i) < 1$  and  $0 < F_N(y_j) < 1$  for all scores  $x_i$  (defaulters) and  $y_j$  (survivors) in the estimation sample. This will not be the case if  $F_N$  is estimated as the empirical distribution function of the survivor sample  $y_j$ . To work around this issue, the empirical distribution can be modified (as described in section 5.3). Another approach could be non-linear least squares estimation:

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta \in \mathbb{R}} \mathbb{E} \left[ \left( \mathbf{1}_D - \{1 + \exp(\alpha + \beta \Phi^{-1}(F_N(S)))\}^{-1} \right)^2 \right], \quad (5.15)$$

where  $\mathbf{1}_D = 1$  for defaulted borrowers and  $\mathbf{1}_D = 0$  otherwise.

Van der Burgt (2008, equation (5)) suggested inferring the value of parameter  $\kappa$  specifying his CAP curve approximation (5.7a) from an estimate of AUC. This idea can be used to derive another approach to the estimation of the parameters in (5.13b) or (5.14). To infer the values of two parameters, two equations are needed. A natural choice for the first equation is to equate a target value  $q$  for the unconditional PD and the mean of the conditional PDs:

$$q = \mathbb{E}[\mathbb{P}[D | S]]. \quad (5.16a)$$

This equation can in general be used for the calibration of conditional PDs, see appendix A for details. A good choice for the second equation seems equating a target value  $A$  for the area under the curve  $\text{AUC}^*$  and a representation of  $\text{AUC}^*$  in terms of the conditional PDs:

$$A = \frac{\mathbb{E} \left[ \left( \mathbb{E}[\mathbb{P}[D | S] \mathbf{1}_{\{S < s\}}] \Big|_{s=S} + \mathbb{P}[S = s] \Big|_{\{s=S\}} \mathbb{P}[D | S] / 2 \right) (1 - \mathbb{P}[D | S]) \right]}{\mathbb{E}[\mathbb{P}[D | S]] (1 - \mathbb{E}[\mathbb{P}[D | S]])} \quad (5.16b)$$

This representation of  $\text{AUC}^*$  follows from proposition 3.15. Combining equations (5.16a) and (5.16b) for the inference of parameters can be regarded as a *quasi moment matching* approach. It is “quasi” moment matching because  $\text{AUC}^*$  is not a proper moment of the conditional PDs. The most natural alternative, the variance of the conditional PDs, however, depends on the proportion of defaulters in the borrower population. As this proportion clearly varies over time it would be difficult to determine an appropriate target variance of the conditional PDs. In contrast,  $\text{AUC}^*$  by its definition does not depend on the proportion of defaulters in the borrower population. It is therefore plausible to assume that discriminatory power displays less variation over time such that its value can be inferred from a historical estimation sample and still applies to the current portfolio. The following example illustrates how the quasi moment matching approach works when the logit shape (5.14) is supposed to apply to the conditional PDs.

**Example 5.4 (Quasi moment matching for PD curve calibration)**

Let  $s_1 \leq s_2 \leq \dots \leq s_n$  be a sorted calibration sample of credit scores in the sense of section 2.1.2 (possibly the scores of the current portfolio). Assume that the PDs conditional on the score

realisations can be described by (5.14). The sample versions of equations (5.16a) and (5.16b) then read:

$$\begin{aligned}
 q &= \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + \exp(\alpha + \beta s_i)}, \\
 A &= \frac{\sum_{i=1}^n \frac{\exp(\alpha + \beta s_i)}{1 + \exp(\alpha + \beta s_i)} \left( \frac{1}{2(1 + \exp(\alpha + \beta s_i))} + \sum_{j=1}^{i-1} \frac{1}{1 + \exp(\alpha + \beta s_j)} \right)}{\left( \sum_{i=1}^n \frac{1}{1 + \exp(\alpha + \beta s_i)} \right) \left( \sum_{i=1}^n \frac{\exp(\alpha + \beta s_i)}{1 + \exp(\alpha + \beta s_i)} \right)}.
 \end{aligned} \tag{5.17}$$

Here  $q$  is the target unconditional PD which could be estimated for instance by econometric methods (see, e.g., Engelmann and Porath, 2003, section III). The variable  $A$  stands for the target discriminatory power, expressed as area under the curve AUC\* which can be estimated from an estimation sample in the sense of section 2.1.1 by means of (3.36a).

Solving the equation system (5.17) for the parameters  $\alpha$  and  $\beta$  then gives the quasi moment matching coefficients for the logit approach to conditional PDs.

### 5.3 Performance comparison

To illustrate the operation of van der Burgt's approach and the three logit approaches introduced in section 5.2 and to compare their performance, we get back to the example from section 4.2 for the continuous score distribution case and to the example from section 4.3 for the case of a discrete rating distribution. The examples, together with some modifications, will show that none of the four approaches is uniformly superior to the others. To see how the estimation methods work we conduct simulation experiments with five different scenarios.

The following scenarios are considered:

1) Rating systems with discrete scales:

- Case 1: 17 grades, binomial distribution with probability parameter 0.4 for the defaulters' rating distribution, binomial distribution with probability parameter 0.5 for the survivors' rating distribution (as in section 4.3).
- Case 2: 7 grades, binomial distribution with probability parameter 0.3 for the defaulters' rating distribution, binomial distribution with probability parameter 0.5 for the survivors' rating distribution.

2) Score functions with continuous scales:

- Case 3: Normal distribution with mean 6.8 and standard deviation 1.96 for the defaulters' score distribution, normal distribution with mean 8.5 and standard deviation 2 for the survivors' score distribution (as in section 4.2). Means and standard deviations are chosen such as to match those from the above discrete case 1.
- Case 4: Normal distribution with mean 2.1 and standard deviation 1.12 for the defaulters' score distribution, normal distribution with mean 3.5 and standard deviation 1.22 for the survivors' score distribution. Means and standard deviations are chosen such as to match those from the above discrete case 2.
- Case 5: Normal distribution with mean 0.0 and standard deviation 1.25 for the defaulters' score distribution, normal distribution with mean 1.0 and standard deviation

1.0 for the survivors' score distribution. Means and standard deviations here are chosen such as to have a larger difference in standard deviations than in cases 3 and 4 and to have the standard deviation for the defaulters' score distribution exceeding the standard deviation for the survivors' score distribution.

In cases 1 and 2, when the value of the unconditional PD is known, the true conditional PDs per rating grade can be calculated according to (2.6a). In cases 3, 4, and 5 the true conditional PDs per score value can be calculated according to (2.6b).

For each scenario the following simulation experiment with 1000 iterations is conducted:

- 1) Generate an estimation sample: Rating grades / scores of 25 (results in table 7) and 50 (results in table 8) defaulters and rating grades / scores of 250 survivors.
- 2) Based on the estimation sample, estimates are calculated for
  - discriminatory power measured by  $AUC^*$  according to (3.36a),
  - parameters<sup>12</sup>  $p$  and  $\kappa$  and distribution function<sup>13</sup>  $F$  for the raw conditional PDs suggested by van der Burgt (2008), where the PDs are calculated according to (5.9b) in the continuous cases 3, 4, and 5, and according to (5.10b) in the discrete cases 1 and 2,
  - parameters  $\alpha$  and  $\beta$  and distribution function<sup>14</sup>  $F_N$  for the raw conditional PDs according to the robust logit approach (5.13b),
  - parameters<sup>15</sup>  $\alpha$  and  $\beta$  for the raw conditional PDs according to the logit approach (5.14).
- 3) Generate then a calibration sample: Rating grades / scores of 300 borrowers with unknown future solvency states. For each of the borrowers first a default / survival simulation with  $PD = 2.5\%$  is conducted. According to the result then a rating grade / score is drawn from the corresponding rating / score distribution. The simulated solvency state is not recorded.
- 4) Based on the calibration sample, the raw PDs from step 2) are calibrated to an unconditional  $PD^{16}$  of 2.5%, as described in proposition A.1 from appendix A.
- 5) Based on the calibration sample, parameters  $\alpha$  and  $\beta$  for PDs according to the quasi moment matching approach are inferred from an unconditional PD 2.5% and the  $AUC^*$  estimate from step 2), as described in example 5.17.

---

<sup>12</sup>Here  $p$  is actually a constant:  $p = 25/(25 + 250) = 1/11$  and  $p = 50/(25 + 250) = 1/6$  respectively.

<sup>13</sup>In the continuous cases 3, 4, and 5, an estimate of  $F$  is calculated according to (2.4).  $F_D$  and  $F_N$  in (2.4) are calculated from normal kernel density estimates with bias-correction as described in section 3.1.2. In the discrete cases 1 and 2 the standard empirical distribution function is deployed for estimating  $F$ .

<sup>14</sup>In the continuous cases 3, 4, and 5,  $F_N$  is calculated from a normal kernel density estimate with bias-correction as described in section 3.1.2. In the discrete cases 1 and 2,  $F_N$  is estimated as the mean of the right-continuous and the left-continuous versions of the empirical distribution function. Additionally, to avoid numerical issues when deploying maximum likelihood estimation, whenever the result would be zero it is replaced by half of the minimum positive value of the modified empirical distribution function.

<sup>15</sup> $\alpha$  and  $\beta$  are estimated by the standard logit MLE procedure (see, e.g., Cramer, 2003, chapter 3).

<sup>16</sup>Hence, we implicitly assume that we have estimated exactly the true unconditional PD of 2.5%. Of course, in practice this would be unlikely. For the purpose of comparing the performance of different estimators this assumption is nonetheless useful.

- 6) Based on the calibration sample, for each rating grade / score the differences between true conditional PD and the PD estimates according to the four different approaches are calculated.
- 7) Based on the four samples of PD differences, the standard error SE is calculated for each of the four approaches according to the generic formula

$$\text{SE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\text{P}[D | S = s_i] - \widehat{\text{P}}[D | S = s_i])^2}, \quad (5.18)$$

where  $s_1, \dots, s_n$  is the calibration sample,  $\text{P}[D | S = s_i]$  is the true conditional PD, and  $\widehat{\text{P}}[D | S = s_i]$  is the estimated conditional PD.

Actually, running the simulation as described in steps 1) to 7) only once would provide interesting insights but would not help too much when it comes to a comparison of the performances of the different estimators. For illustration, have a look at figure 9 which displays the results (true conditional PDs and estimated conditional PDs) of one simulation run of case 1 (17 rating grades scenario). All four estimates seem to fit well the true conditional PDs for rating grades 5 to 17. For rating grades 1 to 4 the fit seems much worse. The van der Burgt and the robust logit estimators even assign constant conditional PD estimates to rating grades 1 to 3.

Note, however, that in this simulation run there were no observations of rating grades 1 to 3 and 16 and 17. In so far, on the one hand, it is questionable whether there should be at all any conditional PD estimates for grades 1 to 3 and 16 and 17. On the other hand, it is not surprising that also for the logit and quasi moment matching estimators the fit at grades 1 to 3 and 16 and 17 is rather poor. Given the sizes of 300 or less of the estimation and calibration samples that are simulated, full coverage of the rating scale by realised rating grades can only be expected for a rating scale with a significantly lower number of grades. In the following, therefore, we look also at the scenario case 2 – a rating system with 7 grades only. The probability to observe unoccupied rating grades when the sample size is about 250 is quite low under scenario case 2.

According to the single simulation run underlying figure 9 there might not be any dramatic differences of the performances of the different estimators. More detailed information can be obtained from running a number of simulations. Tables 7 and 8 (for defaulter scores sample sizes 25 and 50 respectively in the estimation sample) show quantiles of the distributions of the standard errors according to (5.18) that were observed in 1000 iterations of the experiment.

Observations from tables 7 and 8:

- (i) When comparing the quantiles of the standard error distributions as displayed in table 8 to the results from table 7, it appears that the reductions in the low quantiles are moderate while the reductions in the higher quantiles are significantly larger. This observation indicates that the higher number of defaulter scores in the estimation sample mainly has an impact on the variance of the standard error distributions. Note that the distributional assumptions in cases 1 to 5 have been chosen deliberately such that exact matches by one of the estimation approaches are not possible. Hence the standard error rates do not converge to zero for larger defaulter score samples. Rather the variances of their distributions will be diminished.
- (ii) For cases 1 to 3 the logit estimator is best according to the quantiles observed at levels 75% or lower. In case 4, there is no clear picture. In case 5, the robust logit estimator is

best. In cases 1 to 4, the variance of the standard error distribution of the van der Burgt estimator is clearly the least.

- (iii) The error magnitude in case 5 is much higher than in the other cases. This might be due to the fact that the true conditional PD curve is not monotonous as a consequence of the fact that there is a relatively large difference between the variances of the two conditional score distributions.
- (iv) The van der Burgt estimator is less volatile than the other estimators (exception in case 5) but has also a much higher minimum error level. Actually, case 5 has been defined deliberately with a higher variance of the defaulters score distribution in order to challenge the performance of the van der Burgt estimator. For figure 8 indicates that the van der Burgt estimator can adapt to the case where the survivors score distribution has larger variance than the defaulters score distribution but not necessarily to the opposite case.
- (v) Performance of the quasi moment matching estimator is inferior to the performance of the logit estimator but the difference is not large.

Van der Burgt (2008, section 5) described an approach to exploring the sensitivity of the conditional PD curves estimated by means of estimator (5.10a) with regard to uncertainty in the estimation sample. This approach is based on the potential variation of the “concavity” parameter  $\kappa$ . Observation (v) indicates that an analogous approach can be applied to the logit estimator by exploring the sensitivities of the quasi moment estimates with respect to  $AUC^*$  and the unconditional PD.

**Remark 5.5 (Use of the quasi moment matching estimator for sensitivity analysis)**

*We have seen in section 4 how to construct confidence intervals for the area under the curve (and equivalently for the accuracy ratio) even in case of defaulter scores sample sizes as small as five. By applying the quasi moment matching estimator, we can then generate conditional PD curves from different values for  $AUC^*$  as indicated by the confidence interval. Similarly, one can vary the unconditional PD which is the other input to the quasi moment matching estimator in order to investigate the sensitivity of the conditional PD curves with respect to the unconditional PD.*

Table 9 displays the results of another approach to the performance comparison of the conditional PD estimates. The table shows for both defaulter score sample sizes of 25 and 50 and all the five scenarios introduced earlier the frequencies (in 1000 simulation iterations) with which the four estimation approaches produced the least standard error. Hence the entries for the different estimators in a row of table 9 add up to 100%. The results from table 9 re-affirm observations (ii) and (v) made on tables 7 and 8 in so far as they also show dominance of the logit or quasi moment matching estimators in cases 1 to 3 and of the robust logit estimator in case 5. Table 9 however, indicates a clear superiority of the van der Burgt estimator in case 4 where the results of tables 7 and 8 are less clear.

Also shown in table 9 (last column) are the ratios of the conditional score distribution standard deviations for the five considered scenarios. This helps to explain the performance results.

- The logit and quasi moment matching estimators stand out in cases 1 and 3 because then the standard deviations are nearly equal and therefore (5.14) describes an almost

exact fit of the conditional PD curve (Cramer, 2003, section 6.1). Note from tables 7 and 8 that nevertheless the estimation error realised with a logit or quasi moment matching estimator can be quite large. This can be explained with a sensitivity analysis as described in remark 5.5. Figure 10 illustrates that matching a wrong AUC-specification has quite a dramatic impact on the shape of the estimated conditional PD curve, with the consequence of a high standard error. Although misspecification of the target unconditional PD has a much weaker impact, this observation clearly reveals significant vulnerability of parametric approaches to conditional PD curve estimation by their dependence on assumptions on the shape of the conditional PD curve.

- The van der Burgt estimator stands out in case 4 because it adapts best to a situation where the survivor score variance is significantly larger than the defaulter score variance. See figure 8 for a graphical demonstration of this adaptation property.
- With a view on case 4, it is surprising that the van der Burgt estimator does not stand out in case 2 although the survivor score variance is also significantly larger than the defaulter score variance. The different approaches to the estimation of  $\kappa$  that we apply in cases 2 and 4 – (5.10a) vs. (5.9a) – might explain this observation. Weighted least squares as in (5.9a) presumably comply better with the standard error definition (5.18) which includes implicit weighting similar to (5.9a).
- The robust logit estimator stands out in case 5 because it adapts best to a situation where the survivor score variance is significantly smaller than the defaulter score variance. The robust logit estimator, however, cannot represent non-monotonous conditional PDs either. That is why the fit even by this estimator in case 5 is quite poor (as shown in tables 7 and 8).

As the final observation in this simulation study table 10 shows Spearman rank correlations between the absolute errors of the AUC\* estimates on the estimation sample and the standard errors of the conditional PD estimates on the calibration sample. Again the last column of the table displays the ratios of the conditional score distribution standard deviations for the five considered scenarios. Table 10 demonstrates that there is a clear relation between the two estimation errors if the variances of the conditional score distributions are approximately equal. The less equal the variances of the conditional score distributions are, the weaker the relation seems to be. However, the almost vanishing correlations in case 5 could also be caused by the rather high estimation errors observed for the conditional PDs in this case. Hence, it seems premature to draw a firm conclusion from this limited evidence.

## 6 Conclusions

In this paper, we have treated some topics that are not very closely related at first glance:

- 1) In section 3 we have looked in detail at the question of how to define and calculate consistently discriminatory power in terms of area under the curve (AUC) and accuracy ratio (AR). We have seen that there are good reasons to base the definitions of AUC and AR on definition 3.13 of modified ROC and CAP curves. Section 3.2.2 provides ready-to-use formulas for the estimation of AUC and AR from samples.

- 2) In section 4 we have reported the results of a simulation study which refined related studies by Engelmann et al. (2003a,b) on the performance of confidence interval estimators for AUC. We have confirmed that the asymptotic normality confidence interval estimator is most reliable. However, not surprisingly even this estimator performs not very well when applied to defaulter score samples of size ten or less.
- 3) In section 5 we have discussed a proposal by van der Burgt (2008) to derive PD curve estimates by a one-parameter approach to the estimation of CAP curves. By providing background information, we have shown that there are some quite natural logit-related alternatives to van der Burgt's proposal. We have then investigated the performance of the different estimators by another simulation study. The results of this study are mixed in that they demonstrate on the one hand that none of the discussed estimation methods is uniformly best and on the other hand that, in general, by following a parametric approach one risks huge estimation errors caused by the implicit structural assumptions of the estimators.

The common theme in this list is the fact that all the estimation concepts and techniques can be deployed in an implementation of the two-phases approach to PD curve estimation as described in section 2.1. In the first phase of this approach one estimates shape parameters that are essentially equivalent to discriminatory power as expressed by AUC or AR (van der Burgt's concavity parameter  $\kappa$  or the parameter  $\beta$  in the logit curves from section 5.2). In the second phase of the approach the *raw PD curve* from the first phase is calibrated on the current portfolio such that the resulting unconditional probability of default fits an independently determined target unconditional PD. The technical details of this calibration step are described in appendix A.

## References

- G. Casella and R. L. Berger. *Statistical Inference*. Duxbury Press, second edition, 2002.
- B. Clavero Rasero. *Statistical Aspects of Setting up a Credit Rating System*. PhD thesis, Fachbereich Mathematik, Technische Universität Kaiserslautern, 2006.
- J. S. Cramer. *Logit Models From Economics and Other Fields*. Cambridge University Press, 2003.
- A. C. Davison and D. V. Hinkley. *Bootstrap Methods and their Application*. Cambridge University Press, 1997.
- R. Durrett. *Probability: Theory and Examples*. Duxbury Press, second edition, 1995.
- B. Engelmann and D. Porath. Empirical Comparison of Different Methods for Default Probability Estimation. Working paper, 2003.
- B. Engelmann and R. Rauhmeier, editors. *The Basel II Risk Parameters: Estimation, Validation, and Stress Testing*. Springer, 2006.
- B. Engelmann, E. Hayden, and D. Tasche. Testing rating accuracy. *RISK*, 16(1):82–86, January 2003a.

- B. Engelmann, E. Hayden, and D. Tasche. Measuring the Discriminative Power of Rating Systems. Discussion paper (Series 2: Banking and Financial Supervision) 01/2003, Deutsche Bundesbank, 2003b.
- T. Fawcett. ROC Graphs: Notes and Practical Considerations for Researchers. Working paper, 2004. URL [http://home.comcast.net/~tom.fawcett/public\\_html/papers/ROC101.pdf](http://home.comcast.net/~tom.fawcett/public_html/papers/ROC101.pdf).
- D. J. Hand. *Construction and Assessment of Classification Rules*. John Wiley & Sons, Chichester, 1997.
- R. Newson. Parameters behind “non-parametric” statistics: Kendall’s  $\tau_a$ , somers’  $D$  and median differences. *The Stata Journal*, 1(1):1–20, 2001.
- R. Newson. Confidence intervals for rank statistics: Somers’  $D$  and extensions. *The Stata Journal*, 6(3):309–334, 2006.
- OeNB. *Guidelines on Credit Risk Management: Rating Models and Validation*. Oesterreichische Nationalbank and Austrian Financial Market Authority, November 2004. URL [http://www.oenb.at/en/img/rating\\_models\\_tcm16-22933.pdf](http://www.oenb.at/en/img/rating_models_tcm16-22933.pdf).
- A. Pagan and A. Ullah. *Nonparametric econometrics*. Cambridge University Press, 1999.
- K. Pluto and D. Tasche. Thinking positively. *RISK*, 18:72–78, 2005.
- R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2008. URL <http://www.R-project.org>. ISBN 3-900051-07-0.
- S. Satchell and W. Xia. Analytic models of the ROC curve: Applications to credit rating model validation. In G. Christodoulakis and S. Satchell, editors, *The Analytics of Risk Model Validation*, pages 113–133. Academic Press, 2008.
- D. Tasche. Validation of internal rating systems and PD estimates. In G. Christodoulakis and S. Satchell, editors, *The Analytics of Risk Model Validation*, pages 169–196. Academic Press, 2008.
- M. van der Burgt. Calibrating low-default portfolios, using the cumulative accuracy profile. *Journal of Risk Model Validation*, 1(4):17–33, 2008.
- A. W. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 1998.
- E. Weisstein. Fisher’s Exact Test, 2009. URL <http://mathworld.wolfram.com/FishersExactTest.html>.

## A Appendix: Calibration of a PD curve to a target unconditional PD

We assume that for each score or rating grade  $s$  an estimate  $p_{D,p}(s) = P_p[D | S = s]$  of the PD conditional on the score has been made. The index  $p$  indicates that these PDs depend on the unconditional PD  $p = \frac{n_D}{n_D + n_N}$  where  $n_D$  is the size of the defaulter estimation sample  $x_1, \dots, x_{n_D}$



and  $n_N$  stands for the size of the survivor estimation sample  $y_1, \dots, y_{n_N}$ . The PDs  $p_{D,p}(s)$  are called *raw PDs*. In section 5 we look at some parametric approaches to the estimation of such raw PDs.

It is, however, unlikely that the unknown defaulter proportion (the actual unconditional PD) in a given calibration sample  $s_1, \dots, s_n$  (possibly the current portfolio) is  $p$ . We assume that instead there is an estimate  $\pi \neq p$  of this unknown unconditional PD. The aim is then to find a transformation of the raw PDs evaluated on the sample  $s_1, \dots, s_n$  such that their mean equals  $\pi$ . By (2.5) this is a necessary condition for having unbiased estimates of the conditional PDs.

In the special cases of a rating system with a fixed number of grades  $k$  (i.e.  $S$  is a discrete random variable) and of a score function with conditional score densities  $f_N$  and  $f_D$ , we know from equations (2.6a) and (2.6b) that the unconditional PD  $q$  and the corresponding conditional PDs  $p_{D,q}(s)$  satisfy the following equation:

$$\lambda(s) = \frac{q}{1-q} \frac{1 - p_{D,q}(s)}{p_{D,q}(s)}, \quad (\text{A.1})$$

where the *likelihood ratio*  $\lambda$  is defined as

$$\lambda(s) = \begin{cases} \frac{f_N(s)}{f_D(s)}, & \text{S continuous,} \\ \frac{\text{P}[S = s | N]}{\text{P}[S = s | D]}, & \text{S discrete.} \end{cases}$$

As mentioned in section 2.1.2, we assume that the conditional score distributions are the same in the estimation and in the calibration sample. Then also the likelihood ratios are the same in the estimation and in the calibration sample. Hence (A.1) applies both to the raw PDs with unconditional PD  $p$  and to the conditional PDs  $p_{D,\pi}(s)$  corresponding to the unconditional PD  $\pi$ . This observation implies<sup>17</sup>

$$\begin{aligned} p_{D,\pi}(s) &= \frac{1}{1 + \frac{1-\pi}{\pi} \lambda(s)} \\ &= \frac{1}{1 + \frac{1-\pi}{\pi} \frac{p}{1-p} \frac{1 - p_{D,p}(s)}{p_{D,p}(s)}}. \end{aligned} \quad (\text{A.2})$$

The PDs from (A.2) often will not have the required property that their mean equals  $\pi$ , even if the conditional score distributions for the estimation and the calibration samples are really the same:

$$\pi \neq \frac{1}{n} \sum_{i=1}^n p_{D,\pi}(s_i). \quad (\text{A.3})$$

This is due to the facts

- that  $\pi$  is unlikely to be an *exact* forecast of the unconditional default rate and
- that the sample  $s_1, \dots, s_n$  is the result of randomly sampling from a mixture of the unconditional score distributions. Hence the empirical distribution of the sample is likely to be somewhat different to the theoretical unconditional score distribution as presented in (2.4).

---

<sup>17</sup>This approach seems to be common knowledge (see, e.g., OeNB, 2004, section 5.3).

Depending upon how much different are the conditional score distributions underlying the estimation sample and the calibration sample respectively and how good a forecast for the true unconditional PD the estimate  $\pi$  is, the difference between the left-hand and the right-hand sides of (A.3) can be of quite different magnitudes. It can become quite large in particular if  $\pi$  is not a point-in-time forecast of the unconditional PD but rather an estimate of a through-the-cycle central tendency.

Whatever the magnitude of the difference is, it may be desirable to obtain equality of the both sides of (A.3) by adjusting the conditional PDs on its right-hand side. The obvious approach to this would be to apply a constant multiplier to each of the  $p_{D,\pi}(s_i)$ . This approach, however, on the one hand lacks a theoretical foundation and, on the other hand, has the disadvantage that conditional PD values higher than 100% may be the consequence of multiplication with a constant factor that is possibly greater than 100%.

In a more sophisticated approach the  $p_{D,\pi}(s_i)$  on the right-hand side of (A.3) are replaced by  $p_{D,q}(s_i)$  where  $q$  is chosen in such a way as to match the left-hand side of (A.3) with its right-hand side. In this approach  $p_{D,q}(s_i)$  is specified by (A.2) (with  $\pi$  substituted by  $q$  and  $s$  substituted by  $s_i$ ). Hence  $q$  is a solution of the equation

$$\pi = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + \frac{1-q}{q} \frac{p}{1-p} \frac{1-p_{D,p}(s_i)}{p_{D,p}(s_i)}}. \quad (\text{A.4a})$$

Recall that the raw PDs  $p_{D,p}(s_i)$  are assumed to be known. It is not difficult to see that the actual value of  $p$  in the fraction  $\frac{p}{1-p}$  in (A.4a) does not matter for the values of the transformed PDs because the transformed PDs depend on  $q$  and  $p$  only through the term  $\frac{1-q}{q} \frac{p}{1-p}$ . Hence it is sufficient to consider the simplified (case  $p = 1/2$  in the fraction  $\frac{p}{1-p}$ ) equation

$$\pi = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + \frac{1-q}{q} \frac{1-p_{D,p}(s_i)}{p_{D,p}(s_i)}}. \quad (\text{A.4b})$$

The right-hand side of this equation is continuous and strictly increasing in  $q$  and converges toward 0 for  $q \rightarrow 0$  and toward 1 for  $q \rightarrow 1$ . Therefore there is a unique solution  $q$  for (A.4b).

**Proposition A.1** *Let  $s_1, \dots, s_n$  be a sample of scores or rating grades. Assume that an estimate  $\pi \in (0, 1)$  of the unconditional PD in  $s_1, \dots, s_n$  is given and that there is a set of raw conditional PDs  $p_D(s_1), \dots, p_D(s_n)$  associated with  $s_1, \dots, s_n$ . If at least one of the  $p_D(s_i)$  is greater than 0 and less than 1, then there is a unique solution  $q = q(\pi) \in (0, 1)$  to equation (A.4b). The numbers*

$$\pi_i = \frac{1}{1 + \frac{1-q(\pi)}{q(\pi)} \frac{1-p_D(s_i)}{p_D(s_i)}}, \quad i = 1, \dots, n \quad (\text{A.5})$$

*are called  $\pi$ -calibrated conditional PDs associated with the sample  $s_1, \dots, s_n$ .*

Figure 1: Score densities and ROC and CAP curves in the case of normal conditional score densities (see section 3.1.1). Parameter values as in (4.2). Unconditional PD 10%.

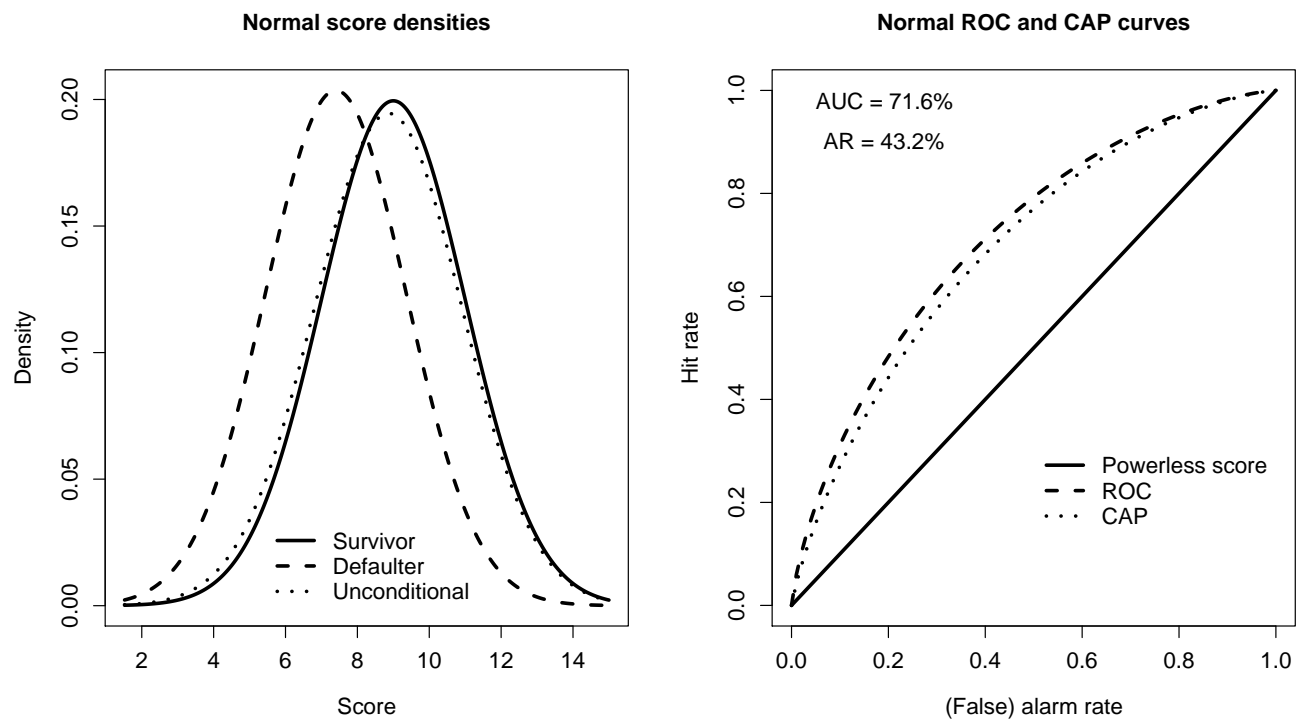


Figure 2: *Non-parametric estimates (with normal kernels) of conditional score densities and ROC curve. Samples of size  $n_D = 5$  and  $n_N = 250$  from normal densities as in figure 1.*

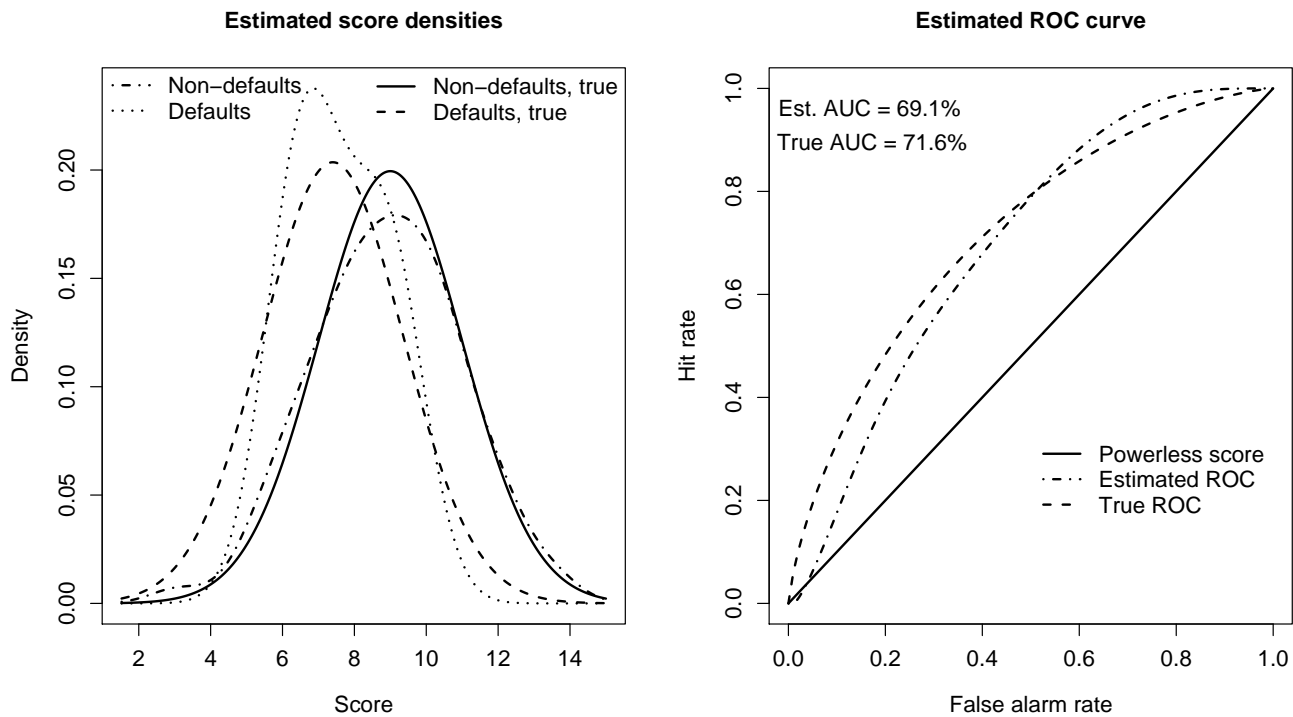


Figure 3: *Fictitious conditional rating distributions for a rating system with 17 grades.*  
*Upper panel: Defaulters' distribution is binomial with success probability 40%; survivors' distribution is binomial with success probability 50%.*  
*Lower panel: Defaulters' distribution by sampling 5 times from defaulters' distribution from upper panel. Survivors' distribution by sampling 250 times from survivors' distribution from upper panel.*  
*Note the different scaling of the y-axis in the two panels.*

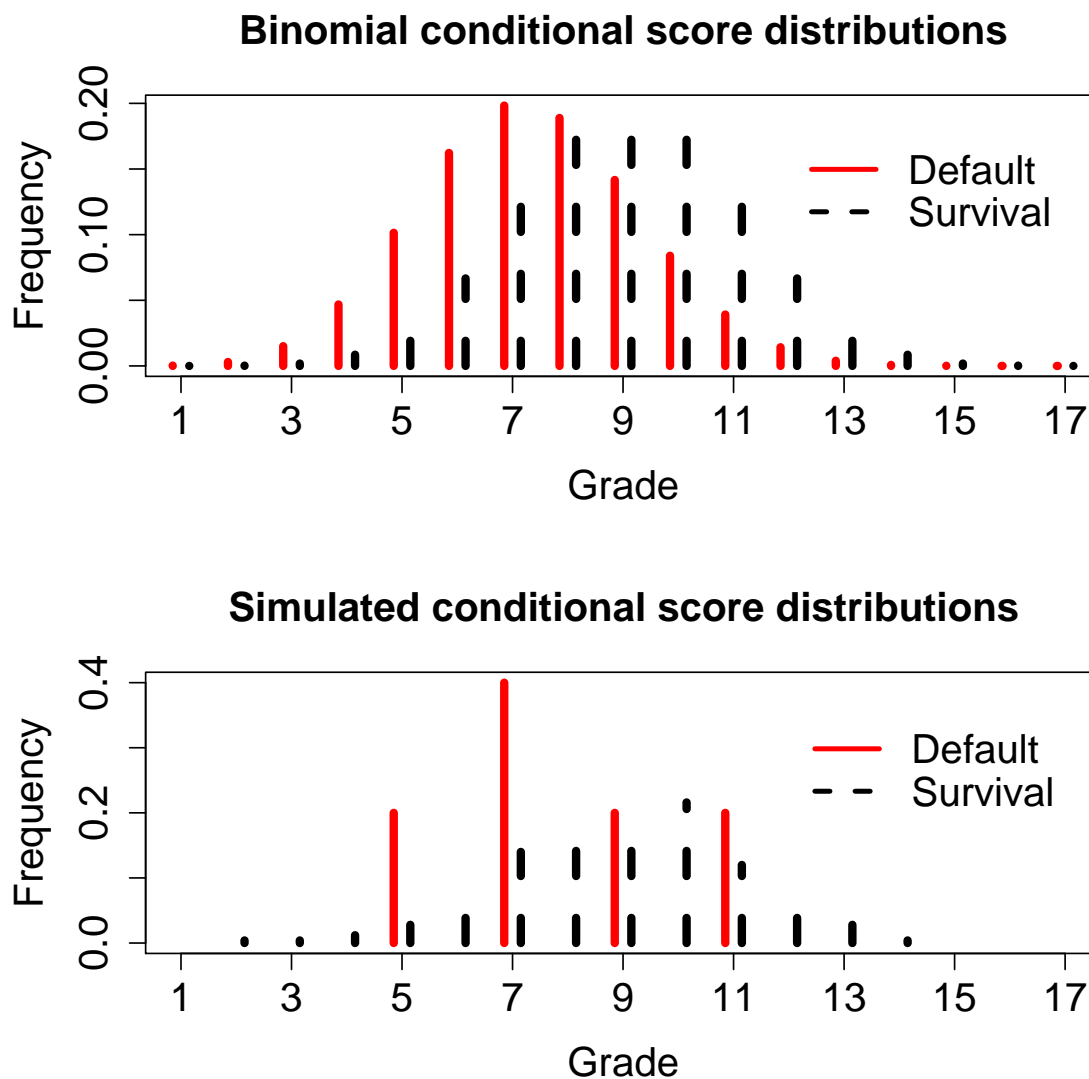


Figure 4: Discrete and modified ROC curves. Conditional rating distributions as in upper panel of figure 3.

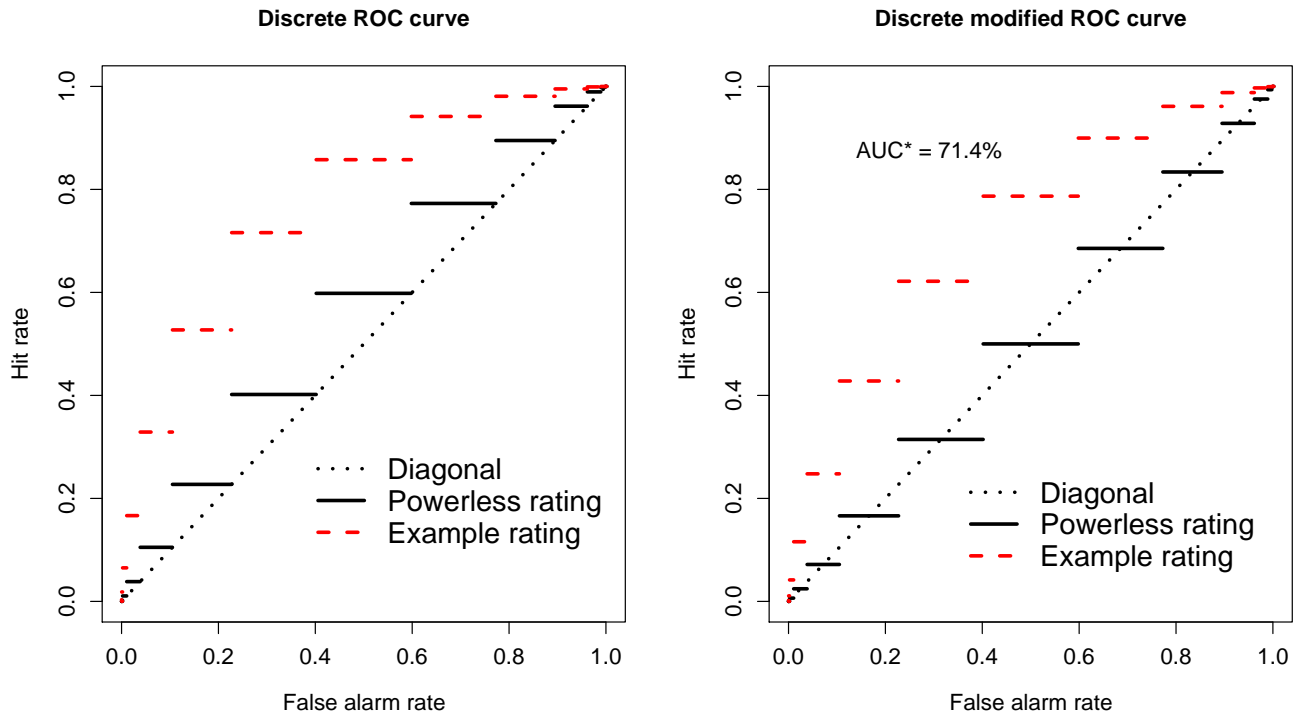


Figure 5: *Modified and interpolated ROC curves. Conditional rating distributions as in lower panel of figure 3.*

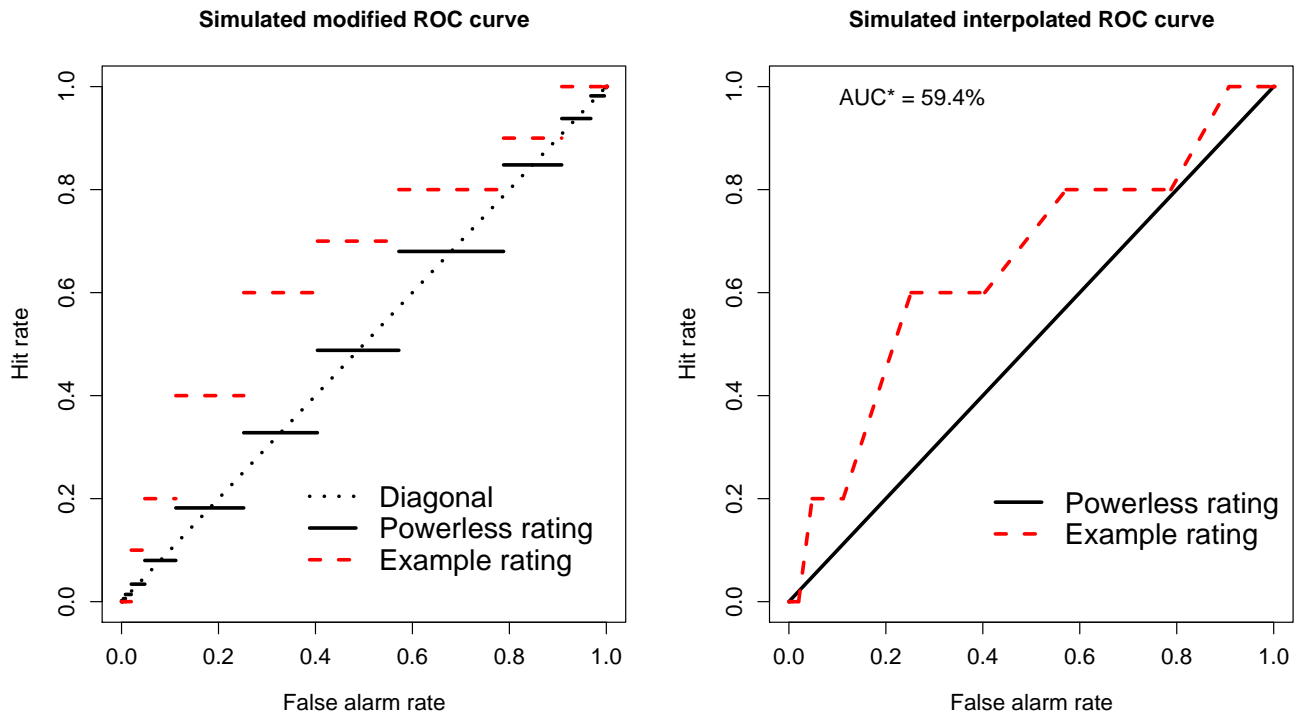


Figure 6: Example from section 4.2. Coverage of true AUC and 50% by 95% confidence intervals as function of sample size  $n_D$  of defaulter scores. Differentiation according to estimation method. Total hits in 100 experiments. Exact results in table 4.

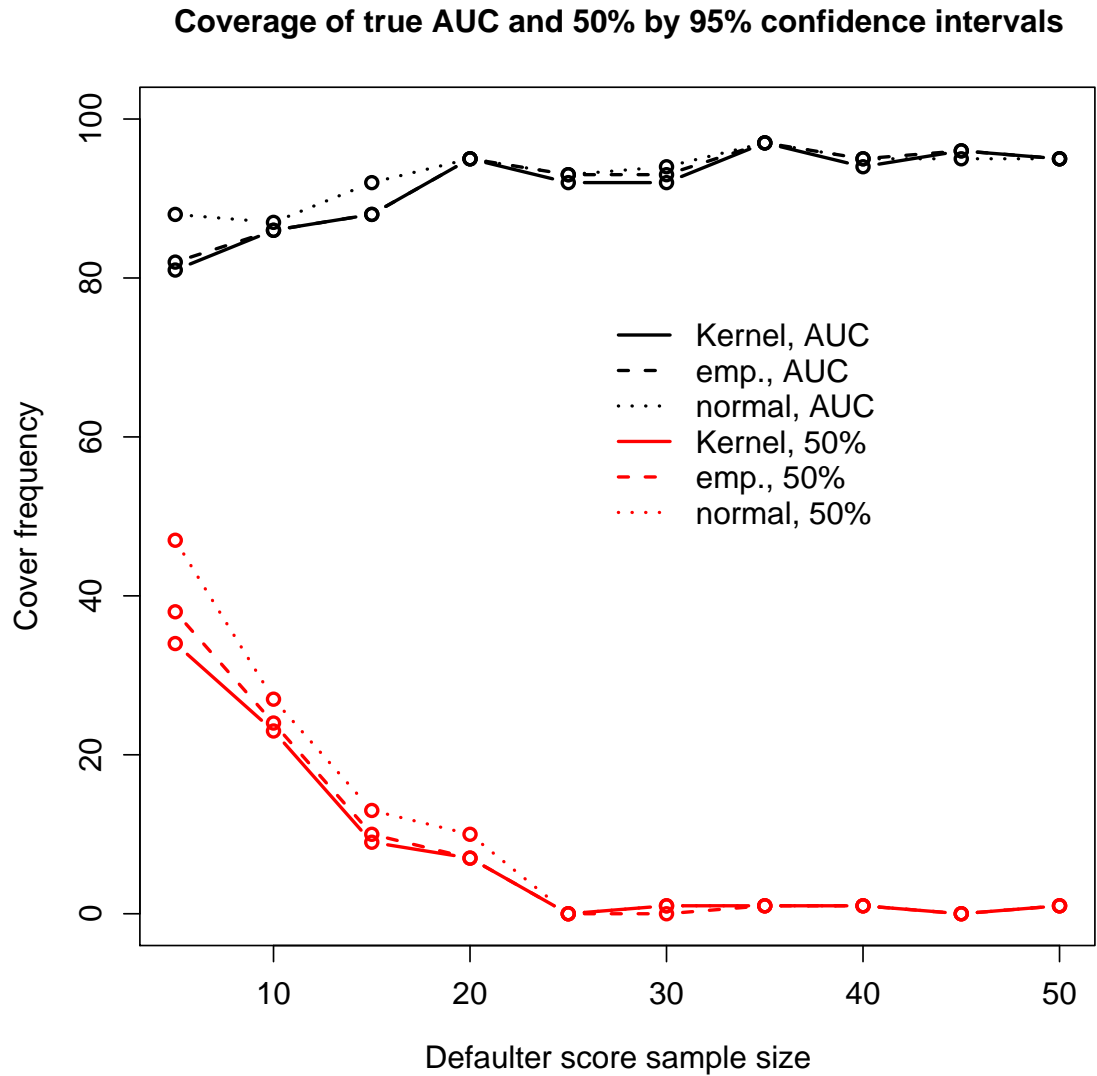




Figure 7: Example from section 4.3. Coverage of true AUC and 50% by 95% confidence intervals as function of sample size  $n_D$  of defaulter scores. Differentiation according to estimation method. Total hits in 100 experiments. Exact results in table 6.

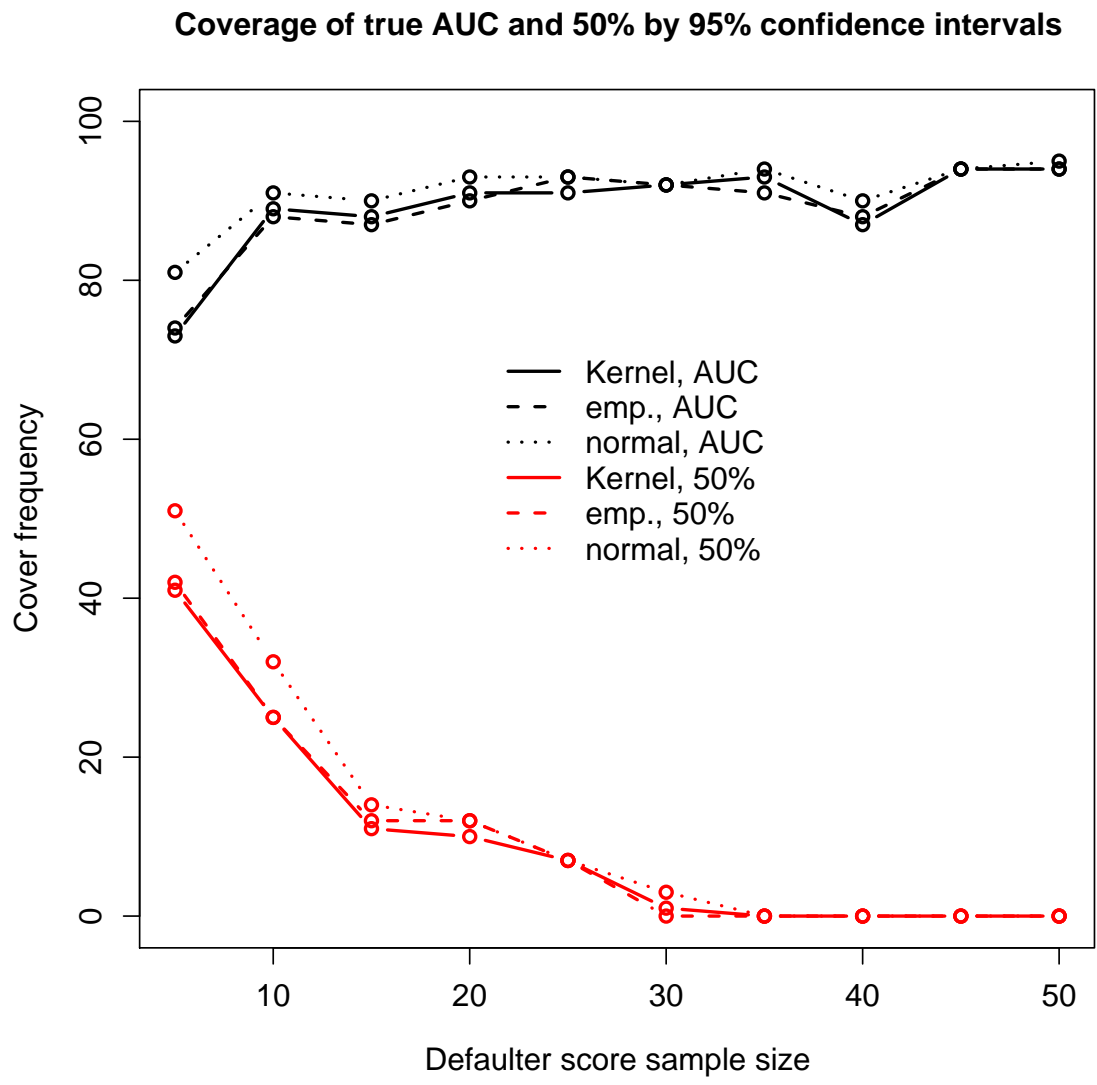


Figure 8: Score densities implied by van der Burgt's parametric approach to CAP curves (5.7a) when the default score distribution is standard normal. The non-default score densities are calculated according to (5.5), with unconditional probability of default  $p = 0.01$ .

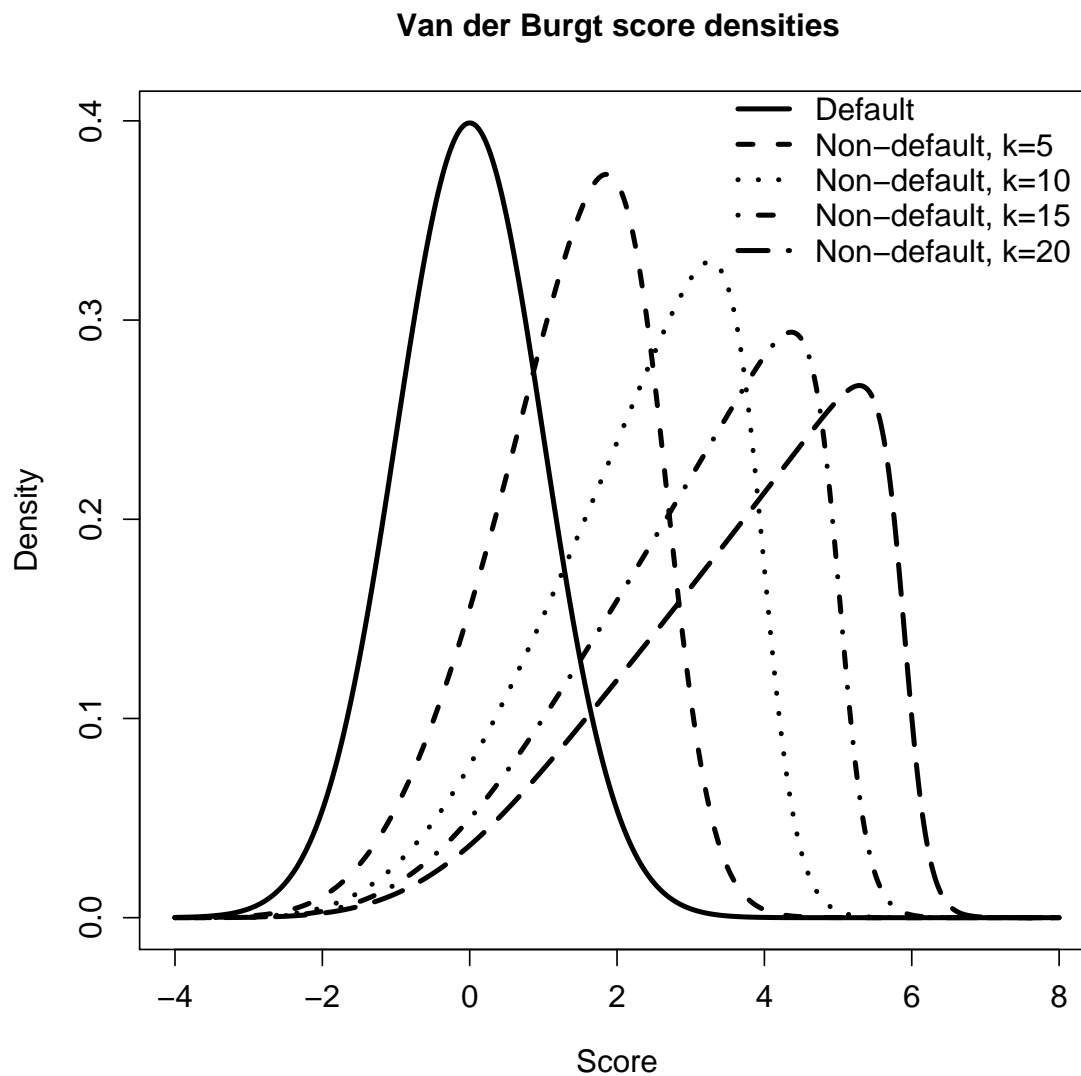


Figure 9: True conditional PDs and estimated conditional PDs for the case 1 scenario (rating system with 17 grades) from section 5.3. Defaulter scores sample size 25 in estimation sample.

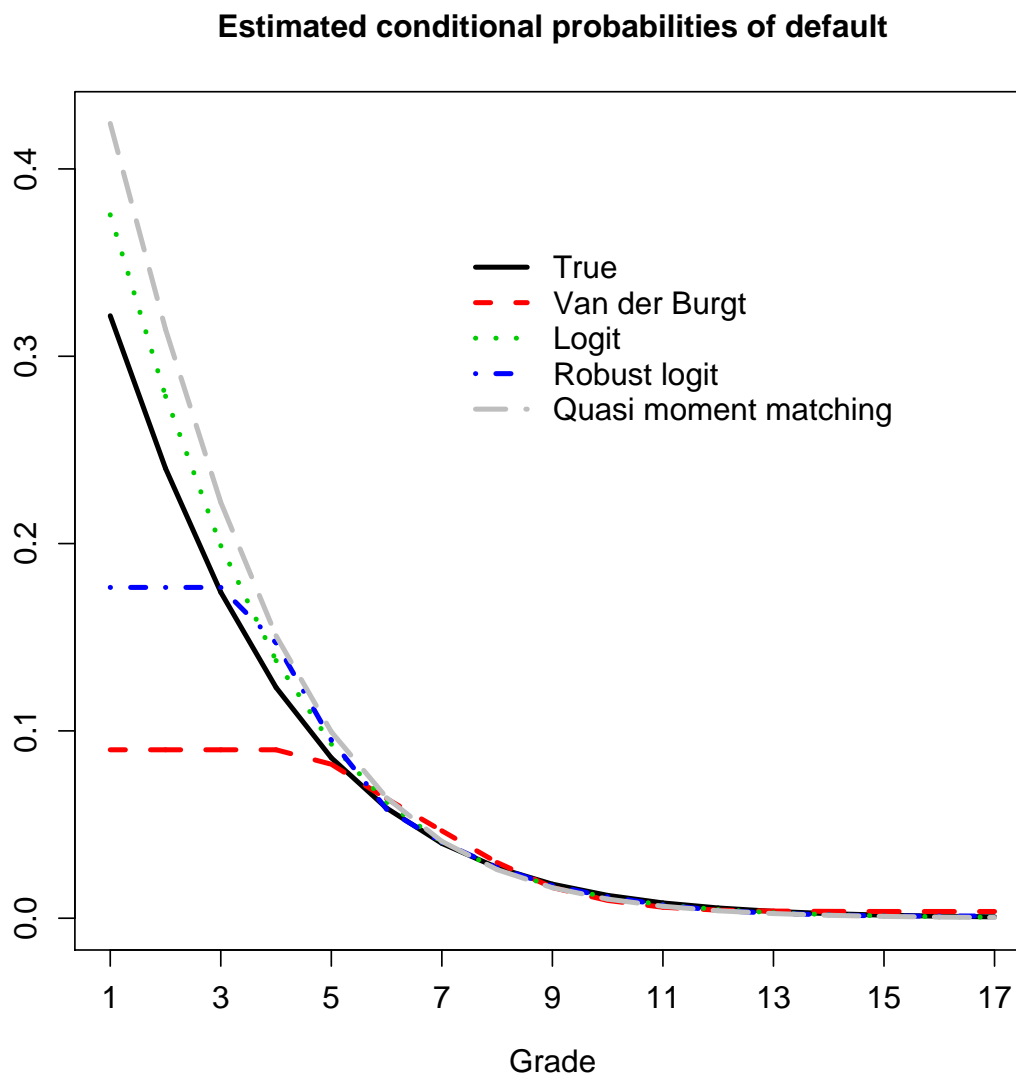


Figure 10: True conditional PDs and conditional PDs estimates by the quasi moment matching estimator (example 5.4) for the case 3 scenario (continuous score function with nearly equal conditional score distribution variances) from section 5.3. Estimates based on calibration sample of size 300. Estimates are matched to the true unconditional PD and the true AUC (“best fit”), to a too small unconditional PD and the true AUC (“lower PD”), and to the true unconditional PD and a too high AUC (“higher AUC”).

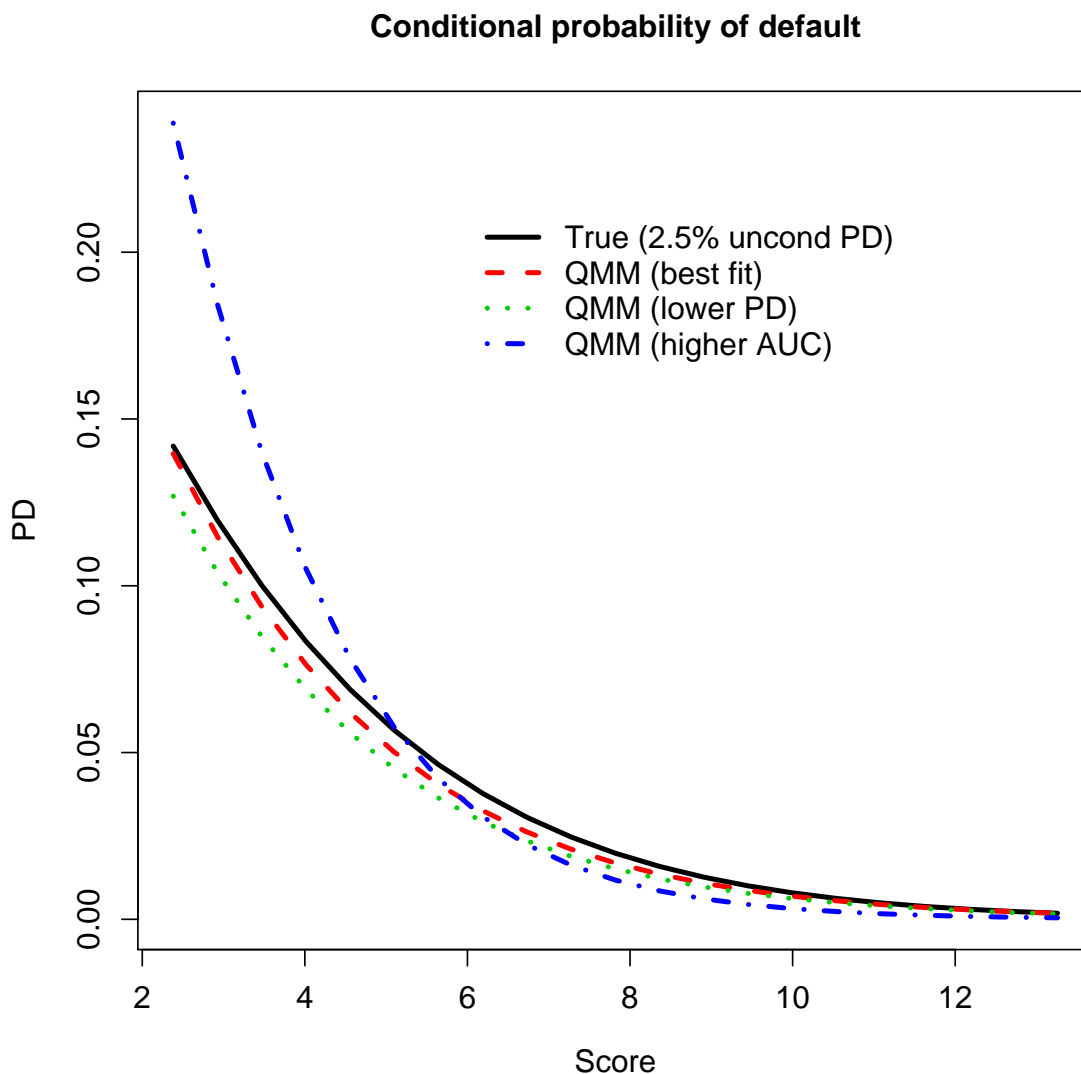


Table 1: Maximum number of different bootstrap samples as function of size  $n$  of original sample.

$n$	1	2	3	4	5	6	7	8	9	10	11
$\binom{2n-1}{n}$	1	3	10	35	126	462	1716	6435	24310	92378	352716

Table 2: Estimated (from 100 simulation experiments) mean numbers  $\mu_n$  and  $\nu_n$  of different (after sorting) samples in 1000 bootstrap iterations of size  $n$  with  $n$  different elements and  $n - 1$  different elements, respectively.

$n$	1	2	3	4	5	6	7	8	9	10	11
$\mu_n$	1.0	3.0	10.0	35.0	117.0	323.0	620.2	844.6	945.8	983.2	995.1
$\nu_n$	NA	1.0	4.0	15.0	52.0	160.1	389.6	679.8	873.0	957.6	987.0

Table 3: Results of first five bootstrap experiments as described in section 4.2. Confidence intervals at 95% level. Default sample sizes  $n_D = 5$ ,  $n_D = 25$ , and  $n_D = 45$ .

True AUC = 71.615%									
Exp. no.	AUC <sub>kernel</sub>	$I_{\text{kernel}}$	AUC <sub>emp</sub>	$I_{\text{emp}}$	$I_{\text{normal}}$				
$n_D = 5$									
1	56.95%	37.19%	80.26%	56.80%	36.96%	81.28%	32.88%	80.72%	
2	68.67%	47.30%	92.75%	68.56%	47.12%	94.08%	42.19%	94.93%	
3	62.14%	42.19%	84.07%	62.32%	42.08%	84.64%	37.52%	87.12%	
4	75.37%	54.99%	97.46%	73.52%	52.00%	95.92%	48.88%	98.16%	
5	62.19%	37.26%	86.07%	63.12%	38.48%	88.96%	34.53%	91.71%	
$n_D = 25$									
1	65.67%	55.21%	77.19%	65.89%	55.25%	77.74%	54.42%	77.36%	
2	65.55%	55.18%	76.22%	65.60%	54.93%	76.85%	54.59%	76.61%	
3	62.44%	51.86%	73.42%	62.18%	51.62%	73.52%	51.04%	73.31%	
4	70.44%	59.02%	81.74%	71.28%	60.91%	82.62%	59.80%	82.76%	
5	66.62%	55.78%	77.62%	66.62%	56.05%	78.27%	55.63%	77.62%	
$n_D = 45$									
1	68.63%	61.30%	77.42%	68.47%	61.12%	77.08%	60.45%	76.49%	
2	72.07%	63.34%	81.71%	71.74%	62.66%	82.15%	62.63%	80.86%	
3	75.00%	68.27%	83.28%	74.95%	68.24%	83.38%	67.35%	82.55%	
4	71.45%	63.67%	79.81%	71.15%	63.04%	79.70%	62.26%	80.03%	
5	67.31%	59.77%	75.60%	67.06%	59.34%	75.52%	59.16%	74.95%	

Table 4: Example from section 4.2. Coverage of true AUC and 50% by 95% confidence intervals. With differentiation according to estimation method and sample size  $n_D$  of defaulter scores. Total hits in 100 experiments. MW means Mann-Whitney test, KS means Kolmogorov-Smirnov test.

Method:	True AUC in interval			50% in interval			Type II error rate	
	Kernel	emp.	normal	Kernel	emp.	normal	MW	KS
$n_D = 5$	81	82	88	34	38	47	57	68
$n_D = 10$	86	86	87	23	24	27	29	39
$n_D = 15$	88	88	92	9	10	13	13	19
$n_D = 20$	95	95	95	7	7	10	10	14
$n_D = 25$	92	93	93	0	0	0	0	10
$n_D = 30$	92	93	94	1	0	1	0	5
$n_D = 35$	97	97	97	1	1	1	1	2
$n_D = 40$	94	95	95	1	1	1	1	0
$n_D = 45$	96	96	95	0	0	0	0	0
$n_D = 50$	95	95	95	1	1	1	1	1

Table 5: Results of first five bootstrap experiments as described in section 4.3. Confidence intervals at 95% level. Default sample sizes  $n_D = 5$ ,  $n_D = 25$ , and  $n_D = 45$ .

True AUC = 71.413%								
Exp. no.	AUC <sub>kernel</sub>	$I_{\text{kernel}}$		AUC <sub>emp</sub>	$I_{\text{emp}}$		$I_{\text{normal}}$	
$n_D = 5$								
1	69.04%	48.88%	94.86%	69.12%	49.64%	95.96%	43.18%	95.06%
2	63.97%	36.87%	92.31%	62.80%	35.04%	92.76%	30.07%	95.53%
3	68.52%	45.21%	96.57%	65.08%	40.00%	89.64%	36.47%	93.69%
4	69.53%	50.58%	92.35%	68.28%	49.16%	89.92%	45.61%	90.95%
5	95.41%	91.81%	100.00%	95.20%	91.40%	100.00%	90.00%	100.00%
$n_D = 25$								
1	68.10%	57.62%	79.41%	68.90%	58.78%	79.82%	57.91%	79.90%
2	69.10%	59.85%	79.31%	68.71%	59.56%	79.29%	58.36%	79.06%
3	69.98%	60.66%	79.51%	69.62%	60.10%	79.43%	59.74%	79.49%
4	66.33%	55.05%	77.66%	66.31%	55.30%	77.65%	54.94%	77.69%
5	80.26%	71.77%	89.89%	79.73%	71.11%	89.89%	70.27%	89.19%
$n_D = 45$								
1	72.00%	64.56%	81.00%	71.64%	64.25%	80.51%	63.75%	79.53%
2	73.34%	66.93%	80.65%	73.03%	66.70%	80.32%	66.20%	79.85%
3	73.22%	65.62%	81.36%	73.14%	65.70%	81.26%	65.36%	80.93%
4	74.13%	66.46%	82.05%	73.92%	66.33%	81.98%	66.17%	81.68%
5	76.63%	70.11%	82.89%	76.32%	69.84%	82.68%	69.74%	82.89%

Table 6: *Example from section 4.3. Coverage of true AUC and 50% by 95% confidence intervals. With differentiation according to estimation method and sample size  $n_D$  of defaulter scores. Total hits in 100 experiments. MW means Mann-Whitney test, Fisher means Fisher's exact test.*

Method:	True AUC in interval			50% in interval			Type II error rate	
	Kernel	emp.	normal	Kernel	emp.	normal	MW	Fisher
$n_D = 5$	73	74	81	41	42	51	63	79
$n_D = 10$	89	88	91	25	25	32	32	61
$n_D = 15$	88	87	90	11	12	14	15	49
$n_D = 20$	91	90	93	10	12	12	10	37
$n_D = 25$	91	93	93	7	7	7	6	30
$n_D = 30$	92	92	92	1	0	3	1	19
$n_D = 35$	93	91	94	0	0	0	0	9
$n_D = 40$	87	88	90	0	0	0	1	8
$n_D = 45$	94	94	94	0	0	0	0	7
$n_D = 50$	94	94	95	0	0	0	0	5

Table 7: *Standard errors according to (5.18) for different approaches to estimation of conditional probabilities of default. Measured in simulation experiment as described in section 5.3. Defaulter scores sample size 25 in estimation sample.*

Quantile level	Standard errors			
	Quasi moment matching	Logit	Robust logit	Van der Burgt
Case 1 (17 rating grades)				
5%	0.097%	0.088%	0.185%	0.483%
25%	0.268%	0.269%	0.339%	0.656%
50%	0.528%	0.496%	0.553%	0.82%
75%	0.9%	0.861%	0.91%	1.062%
95%	1.62%	1.584%	1.7%	1.479%
Case 2 (7 rating grades)				
5%	0.115%	0.116%	0.213%	0.58%
25%	0.32%	0.317%	0.392%	0.717%
50%	0.577%	0.543%	0.657%	0.86%
75%	0.942%	0.915%	1.089%	1.12%
95%	1.628%	1.67%	2.022%	1.593%
Case 3 (continuous ~ 17 grades)				
5%	0.121%	0.123%	0.169%	0.45%
25%	0.285%	0.27%	0.372%	0.621%
50%	0.523%	0.517%	0.644%	0.791%
75%	0.933%	0.908%	1.111%	1.027%
95%	1.75%	1.761%	2.518%	1.512%
Case 4 (continuous ~ 7 grades)				
5%	0.288%	0.28%	0.261%	0.425%
25%	0.471%	0.454%	0.533%	0.599%
50%	0.735%	0.715%	0.886%	0.794%
75%	1.224%	1.185%	1.662%	1.046%
95%	2.235%	2.315%	3.387%	1.559%
Case 5 (continuous, different variances)				
5%	0.646%	0.634%	0.493%	1.133%
25%	0.902%	0.875%	0.782%	1.66%
50%	1.348%	1.24%	1.164%	2.282%
75%	2.18%	1.913%	1.74%	3.238%
95%	3.7%	3.398%	3.028%	4.994%



Table 8: *Standard errors according to (5.18) for different approaches to estimation of conditional probabilities of default. Measured in simulation experiment as described in section 5.3. Defaulter scores sample size 50 in estimation sample.*

Quantile level	Standard errors			
	Quasi moment matching	Logit	Robust logit	Van der Burgt
Case 1 (17 rating grades)				
5%	0.083%	0.085%	0.17%	0.414%
25%	0.216%	0.237%	0.31%	0.561%
50%	0.418%	0.4%	0.483%	0.696%
75%	0.684%	0.697%	0.777%	0.885%
95%	1.168%	1.183%	1.302%	1.243%
Case 2 (7 rating grades)				
5%	0.101%	0.097%	0.182%	0.426%
25%	0.265%	0.267%	0.36%	0.545%
50%	0.477%	0.475%	0.594%	0.68%
75%	0.802%	0.781%	0.943%	0.879%
95%	1.289%	1.288%	1.754%	1.295%
Case 3 (continuous ~ 17 grades)				
5%	0.114%	0.117%	0.163%	0.356%
25%	0.244%	0.245%	0.336%	0.524%
50%	0.437%	0.44%	0.557%	0.676%
75%	0.741%	0.722%	0.928%	0.898%
95%	1.299%	1.275%	2.222%	1.264%
Case 4 (continuous ~ 7 grades)				
5%	0.281%	0.281%	0.259%	0.307%
25%	0.456%	0.426%	0.5%	0.466%
50%	0.695%	0.681%	0.852%	0.643%
75%	1.095%	1.05%	1.458%	0.855%
95%	1.829%	1.92%	3.214%	1.24%
Case 5 (continuous, different variances)				
5%	0.635%	0.634%	0.47%	0.968%
25%	0.872%	0.867%	0.745%	1.485%
50%	1.273%	1.257%	1.082%	2.12%
75%	2.008%	1.923%	1.594%	3.055%
95%	3.351%	3.07%	2.688%	4.638%

Table 9: Probabilities to produce least standard error for different approaches to estimation of conditional probabilities of default. Measured in simulation experiment as described in section 5.3. QMM means “Quasi moment matching”.  $\sigma_D$  and  $\sigma_N$  are the standard deviations of the defaulter score distribution and survivor score distribution respectively.

Case	Probability to produce least standard error				$\sigma_D/\sigma_N$
	QMM	Logit	Robust logit	Van der Burgt	
	Default sample size 25				
1 (17 grades)	40.5%	33.1%	16.6%	9.8%	98%
2 (7 grades)	39%	27.8%	21.7%	11.5%	91.7%
3 (continuous $\sim$ 17 grades)	35.2%	27%	16.9%	20.9%	98%
4 (continuous $\sim$ 7 grades)	29.5%	16.3%	16.7%	37.5%	91.7%
5 (different variances)	25%	21.1%	52.5%	1.4%	125%
	Default sample size 50				
1 (17 grades)	43.7%	31%	14.8%	10.5%	98%
2 (7 grades)	35.8%	28.5%	20.9%	14.8%	91.7%
3 (continuous $\sim$ 17 grades)	35%	28%	14.3%	22.7%	98%
4 (continuous $\sim$ 7 grades)	26%	14.4%	12.9%	46.7%	91.7%
5 (different variances)	25%	17.8%	56%	1.2%	125%

Table 10: Spearman correlations of absolute error of AUC estimate and standard error of conditional PD curve estimate for different approaches to estimation of conditional probabilities of default. Measured in simulation experiment as described in section 5.3. QMM means “Quasi moment matching”.  $\sigma_D$  and  $\sigma_N$  are the standard deviations of the defaulter score distribution and survivor score distribution respectively.

Case	Spearman correlation				$\sigma_D/\sigma_N$
	QMM	Logit	Robust logit	Van der Burgt	
	Default sample size 50				
1 (17 grades)	87.2%	86%	78.1%	55.6%	98%
2 (7 grades)	80.7%	79.5%	65.6%	54.6%	91.7%
3 (continuous $\sim$ 17 grades)	85.4%	85.5%	70.8%	59.3%	98%
4 (continuous $\sim$ 7 grades)	51%	46.9%	38.6%	62.9%	91.7%
5 (different variances)	18.7%	12.4%	25.5%	9.2%	125%
	Default sample size 50				
1 (17 grades)	81.9%	81.8%	71.8%	46.8%	98%
2 (7 grades)	74.4%	75.5%	60.4%	49.6%	91.7%
3 (continuous $\sim$ 17 grades)	79.5%	78.8%	59.7%	49%	98%
4 (continuous $\sim$ 7 grades)	36.4%	30.1%	28.2%	54.8%	91.7%
5 (different variances)	6.9%	4%	11.8%	7.2%	125%